

---

# Large-Scale Sparse Kernel Canonical Correlation Analysis

---

Viivi Uurtio<sup>1,2</sup> Sahely Bhadra<sup>3</sup> Juho Rousu<sup>1,2</sup>

## Abstract

This paper presents gradKCCA, a large-scale sparse non-linear canonical correlation method. Like Kernel Canonical Correlation Analysis (KCCA), our method finds non-linear relations through kernel functions, but it does not rely on a kernel matrix, a known bottleneck for scaling up kernel methods. gradKCCA corresponds to solving KCCA with the additional constraint that the canonical projection directions in the kernel-induced feature space have preimages in the original data space. Firstly, this modification allows us to very efficiently maximize kernel canonical correlation through an alternating projected gradient algorithm working in the original data space. Secondly, we can control the sparsity of the projection directions by constraining the  $\ell_1$  norm of the preimages of the projection directions, facilitating the interpretation of the discovered patterns, which is not available through KCCA. Our empirical experiments demonstrate that gradKCCA outperforms state-of-the-art CCA methods in terms of speed and robustness to noise both in simulated and real-world datasets.

## 1. Introduction

Canonical correlation analysis (CCA) (Hotelling, 1935; 1936) discovers multivariate relations in two-view datasets, by finding linear combinations of variables, the canonical projections, that produce maximum correlation between the two views (Uurtio et al., 2018b).

The standard CCA model has been extended to a non-linear setup by kernel methods (Bach & Jordan, 2002; Hardoon et al., 2004; Chu et al., 2013; Lopez-Paz et al., 2014; Wang

& Livescu, 2016; Uurtio et al., 2018a) and deep learning (Andrew et al., 2013). The advantages of kernelizing the CCA problem, or applying a deep neural network, include the possibility to uncover non-linear multivariate relations. However, the non-linear CCA formulations generally come at the expense of losing the information of which of the variables in the original data matrices are relevant, and typically produce dense models, which hinders interpretability and predictive performance when the underlying relations are sparse. For kernel-based CCA, the quadratic size of the kernel matrix in the number of examples also makes scaling up to big data challenging, although this can be solved through approximations such as random Fourier features (Lopez-Paz et al., 2014) and Nyström approximation (Uurtio et al., 2018a).

This paper proposes a new KCCA variant, named gradKCCA, that is designed for sparse non-linear canonical correlation analysis in large datasets. The contributions of the paper are the following:

- A new model for kernel-based non-linear CCA that maximizes canonical correlation in the kernel-induced feature spaces through the gradients of the preimages of the projection directions, without relying on a kernel matrix. A sparsity-inducing variant of the model is achieved through controlling the  $\ell_1$  norms of the preimages of the projection directions.
- A theoretical analysis of the new model compared to KCCA showing that in the training set, gradKCCA correlation can be bounded between the correlation of KCCA and the correlation achieved by an approximate preimage for KCCA.
- An efficient alternating projected gradient algorithm that has the time-complexity of  $O((p+q)n)$  per epoch, where  $n$  is the number training examples and  $p$  and  $q$  are the numbers of variables in the two views.
- An experimental study including several non-linear CCA variants demonstrating (a) the superior speed and scalability of the method on big data, (b) tolerance to high amounts of irrelevant variables, (c) predictive performance in line with the state-of-the-art on large real world datasets.

<sup>1</sup>Department of Computer Science, Aalto University, Espoo, Finland <sup>2</sup>Helsinki Institute for Information Technology, Helsinki, Finland <sup>3</sup>Computer Science and Engineering, Indian Institute of Technology Palakkad, Palakkad, India. Correspondence to: Viivi Uurtio <viivi.uurtio@aalto.fi>.

## 2. Background

In the following, we briefly review the main approaches proposed for non-linear canonical correlation analysis. These methods will be experimentally compared to our proposed method in Section 5.

**Notation.** Matrices will be denoted by bold upper-case letters (e.g.  $\mathbf{X}$ ,  $\mathbf{Y}$  for data matrices) and vectors by bold lower-case letters.  $\mathbf{I}$  denotes the identity matrix,  $\mathbf{1}$  a is a vector of all ones,  $\|\cdot\|_p$  denotes the  $l_p$  norm, and  $\langle \cdot, \cdot \rangle$  denotes inner product. Multiplication and entry-wise multiplication are denoted by  $\cdot$  and  $\odot$  respectively. We denote sets with non-bold upper-case letters (e.g.  $X$  and  $Y$  for the spaces for data points).

### 2.1. Canonical Correlation Analysis (CCA)

Canonical correlation analysis (Hotelling, 1935; 1936) finds related variables in two-view datasets. Let the two views be given by  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and  $\mathbf{Y} \in \mathbb{R}^{n \times q}$ , where  $n$ ,  $p$ , and  $q$  are the sample size and the numbers of variables in  $\mathbf{X}$  and  $\mathbf{Y}$  respectively. The related variables are determined from the entries of the coefficient vectors  $\mathbf{u} \in \mathbb{R}^p$  and  $\mathbf{v} \in \mathbb{R}^q$  that satisfy

$$\max_{\mathbf{u}, \mathbf{v}} \frac{\langle \mathbf{X}\mathbf{u}, \mathbf{Y}\mathbf{v} \rangle}{\|\mathbf{X}\mathbf{u}\|_2 \|\mathbf{Y}\mathbf{v}\|_2}. \quad (1)$$

In general, the non-convex problem in (1) can be solved through the singular value decomposition (SVD), or in the form of a standard or generalized eigenvalue problem (Urtio et al., 2018b).

The standard techniques work well when the underlying statistical patterns are linear and there are fewer variables than examples. However, it is not possible to model non-linear correlations well using standard CCA.

### 2.2. Kernel Canonical Correlation Analysis (KCCA)

Similarly to the original CCA problem (1), kernel CCA (KCCA) (Bach & Jordan, 2002; Hardoon et al., 2004) is solved through an eigenvalue problem. The  $n$  examples are transformed to Hilbert spaces  $\phi_x : \mathbf{X} \in \mathbb{R}^{n \times p} \mapsto \mathcal{H}_x$  and  $\phi_y : \mathbf{Y} \in \mathbb{R}^{n \times q} \mapsto \mathcal{H}_y$  where  $\phi$  denotes a feature map. The similarities between the observations in the Hilbert spaces are given by symmetric positive semi-definite kernels  $\mathbf{K}_{i,j}^x = \mathbf{K}^x(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi_x(\mathbf{x}_i), \phi_x(\mathbf{x}_j) \rangle_{\mathcal{H}_x}$  and  $\mathbf{K}_{i,j}^y = \mathbf{K}^y(\mathbf{y}_i, \mathbf{y}_j) = \langle \phi_y(\mathbf{y}_i), \phi_y(\mathbf{y}_j) \rangle_{\mathcal{H}_y}$  where  $i, j = 1, 2, \dots, n$ ,  $\mathbf{x} \in \mathbb{R}^p$  for view  $\mathbf{X}$ , and  $\mathbf{y} \in \mathbb{R}^q$  for view  $\mathbf{Y}$  respectively. The kernel canonical correlation is then measured through

$$\max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \rho_{\text{kcca}} = \frac{\langle \mathbf{K}^x \boldsymbol{\alpha}, \mathbf{K}^y \boldsymbol{\beta} \rangle}{\|\mathbf{K}^x \boldsymbol{\alpha}\|_2 \|\mathbf{K}^y \boldsymbol{\beta}\|_2} \quad (2)$$

where  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  are the coefficient vectors to be optimized by KCCA.

The computation of the kernel matrices in KCCA becomes costly at large sample sizes. Matrix decomposition approaches such as the incomplete Cholesky decomposition (Bach & Jordan, 2002; Hardoon et al., 2004) and the partial Gram-Schmidt orthogonalization (PGSO) (Hardoon et al., 2004) have been proposed to speed up the computations. When applying the PGSO in KCCA, every KCCA projection is computed as the span of a subset of the projections of a set of training examples that are selected by performing PGSO of the training vectors in the feature space.

### 2.3. Randomized Non-Linear CCA (RCCA)

The kernel computations of KCCA can also be approximated by random Fourier features (Rahimi & Recht, 2008), as proposed by (Lopez-Paz et al., 2014). In randomized non-linear CCA (RCCA) the kernels are approximated by random Fourier features. For example, as explained in (Lopez-Paz et al., 2014), the examples in view  $\mathbf{X} \in \mathbb{R}^{n \times p} \mapsto \mathbf{z}(\mathbf{X}) \in \mathbb{R}^{n \times m}$  through the following map

$$\begin{aligned} \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m &\sim \mathcal{P}(\mathbf{w}) \\ \mathbf{z}_i &= [\cos(\mathbf{w}_i^\top \mathbf{x}_1 + b_i), \dots, \cos(\mathbf{w}_i^\top \mathbf{x}_n + b_i)] \end{aligned}$$

where  $\mathbf{z} \in \mathbb{R}^n$  for  $i = 1, 2, \dots, m$  and, when the approximated kernel is Gaussian,  $\mathcal{P}(\mathbf{w})$  is a multivariate normal distribution, where  $\mathbf{w} \in \mathbb{R}^d$ .  $b_i$  is sampled from a uniform distribution  $U[0, 2\pi]$ . The  $m$  denotes the number of random Fourier features. The approximated kernel matrix  $\mathbf{K} \approx \hat{\mathbf{K}}$  is given by

$$\hat{\mathbf{K}} = \frac{1}{m} \mathbf{z}(\mathbf{X}) \mathbf{z}(\mathbf{X})^\top = \frac{1}{m} \sum_{i=1}^m \mathbf{z}_i \mathbf{z}_i^\top.$$

The kernel canonical correlations are found through the generalized eigenvalue problem, as in KCCA. In (Lopez-Paz et al., 2014), RCCA was shown to converge to KCCA solution when the number of random Fourier features increases.

### 2.4. Kernel Non-Linear Orthogonal Iterations (KNOI)

The approximation of RCCA to KCCA relies on the dimensionality of the random Fourier feature space. In some cases, a sufficiently good approximation may require a large number of features. This incurs high memory requirements to solve the CCA problem in the random Fourier feature space. A memory-efficient stochastic optimization algorithm of RCCA was proposed in (Wang & Livescu, 2016). The approach was named kernel non-linear orthogonal iterations (KNOI). KNOI employs ideas from the Augmented Approximate Gradient (AppGrad) (Ma et al., 2015) and Nonlinear Orthogonal Iterations (NOI) (Wang et al., 2015) techniques. The idea of KNOI is to update the CCA projections in small mini-batches, using a convex combination of the previous and current estimates. The estimates are then applied to

update the cross-view least squares regression problems. In (Wang & Livescu, 2016), KNOI was shown to outperform RCCA in terms of test correlation and running time, with GPU support.

### 2.5. Deep Canonical Correlation Analysis (DCCA)

Large-scale CCA problems can also be solved using deep neural networks. The deep learning methods apply the standard SVD technique to compute the maximal correlation between the outputs of two neural networks. In deep CCA (DCCA), (Andrew et al., 2013), the  $n$  observations,  $\mathbf{x}_i \in \mathbb{R}^p$  and  $\mathbf{y}_i \in \mathbb{R}^q$  for  $i = 1, 2, \dots, n$ , are transformed several times through a neural network. At every layer, an observation, for example  $\mathbf{x}$ , is transformed by  $\mathbf{h} = s(\mathbf{S}\mathbf{x} + \mathbf{b})$ , where  $\mathbf{S} \in \mathbb{R}^{\text{output} \times \text{input}}$  is a matrix of weights,  $\mathbf{b} \in \mathbb{R}^{\text{output}}$  is a vector of bias, and  $s : \mathbb{R} \mapsto \mathbb{R}$  is generally a non-linear function applied to every entry. The aim is to learn the optimal parameters  $\mathbf{S}$  and  $\mathbf{b}$  for both views such that the correlation between the transformed observations is maximized. Let  $\mathbf{H}_x \in \mathbb{R}^{\text{output} \times n}$  and  $\mathbf{H}_y \in \mathbb{R}^{\text{output} \times n}$  denote the matrices that have the final transformed output vectors in their columns. Let  $\tilde{\mathbf{H}}_x = \mathbf{H}_x - \frac{1}{n}\mathbf{H}_x\mathbf{1}$  denote the centered data matrix and let  $\hat{\mathbf{C}}_{xy} = \frac{1}{m-1}\tilde{\mathbf{H}}_x\tilde{\mathbf{H}}_y^T$ ,  $\hat{\mathbf{C}}_{xx} = \frac{1}{m-1}\tilde{\mathbf{H}}_x\tilde{\mathbf{H}}_x^T + r_x\mathbf{I}$ , and  $\hat{\mathbf{C}}_{yy} = \frac{1}{m-1}\tilde{\mathbf{H}}_y\tilde{\mathbf{H}}_y^T + r_y\mathbf{I}$  where  $r_x$  and  $r_y$  are regularization constants, denote the cross-covariance and within-set covariance matrices. The correlation of the  $k^{\text{th}}$  components is the  $k^{\text{th}}$  singular value of the matrix

$$\mathbf{T} = \hat{\mathbf{C}}_{xx}^{-1/2}\hat{\mathbf{C}}_{xy}\hat{\mathbf{C}}_{yy}^{-1/2}.$$

### 2.6. CCA through Hilbert-Schmidt Independence Criterion (SCCA-HSIC)

Recently, as another direction of research, the Hilbert-Schmidt Independence Criterion (HSIC) (Gretton et al., 2005) has been proposed as a dependence metric for the CCA projections (Chang et al., 2013; Uurtio et al., 2018a). The sparse non-linear variant, named SCCA-HSIC, solves the following optimization problem (Uurtio et al., 2018a):

$$\begin{aligned} \max_{\mathbf{u}, \mathbf{v}} \quad & \frac{\text{trace}(\hat{\mathbf{K}}^{\mathbf{u}}\hat{\mathbf{K}}^{\mathbf{v}})}{(n-1)^2} \\ \text{s.t.} \quad & \|\mathbf{u}\|_1 \leq s_x \\ & \|\mathbf{v}\|_1 \leq s_y \end{aligned} \quad (3)$$

where the values of the associated kernel matrices are given by  $\mathbf{K}_{ij}^{\mathbf{u}} = \langle \mathbf{u}^T \mathbf{x}_i, \mathbf{u}^T \mathbf{x}_j \rangle$  and  $\mathbf{K}_{ij}^{\mathbf{v}} = \langle \mathbf{v}^T \mathbf{y}_i, \mathbf{v}^T \mathbf{y}_j \rangle$ .  $\hat{\mathbf{K}}^{\mathbf{u}}$  and  $\hat{\mathbf{K}}^{\mathbf{v}}$  denote centered kernel matrices. Due to the  $O(n^2)$  size of the kernel matrices, SCCA-HSIC is approximated for large-scale experiments by the Nyström approximation of the kernel matrices (Williams & Seeger, 2001; Zhang et al., 2018).

## 3. gradKCCA: A method for sparse non-linear CCA

In this paper, we consider solving a KCCA problem with the additional constraint that the projection directions  $\mathbf{w}_x \in \mathcal{H}_x$  and  $\mathbf{w}_y \in \mathcal{H}_y$  need to have preimages in the original data spaces, that is,  $\phi_x^{-1}(\mathbf{w}_x) \in \mathbb{R}^p$  and  $\phi_y^{-1}(\mathbf{w}_y) \in \mathbb{R}^q$ . This requirement can alternatively be expressed as requiring the existence of a such  $\mathbf{u} \in \mathbb{R}^p$  that  $\phi_x(\mathbf{u}) = \mathbf{w}_x$  (resp.  $\mathbf{v} \in \mathbb{R}^q$  s.t.  $\phi_y(\mathbf{v}) = \mathbf{w}_y$ ). The resulting CCA problem is written as an optimization over the variables  $\mathbf{u} \in \mathbb{R}^p$ ,  $\mathbf{v} \in \mathbb{R}^q$ :

$$\rho = \max_{\mathbf{u}, \mathbf{v}} \frac{\sum_{i=1}^n \langle \phi_x(\mathbf{x}_i), \phi_x(\mathbf{u}) \rangle \cdot \langle \phi_y(\mathbf{y}_i), \phi_y(\mathbf{v}) \rangle}{\|(\langle \phi_x(\mathbf{x}_i), \phi_x(\mathbf{u}) \rangle)_{i=1}^n\|_2 \|(\langle \phi_y(\mathbf{y}_i), \phi_y(\mathbf{v}) \rangle)_{i=1}^n\|_2} \quad (4)$$

Replacing the inner products with kernel functions  $k^x(\mathbf{x}, \mathbf{z}) = \langle \phi_x(\mathbf{x}), \phi_x(\mathbf{z}) \rangle$ ,  $k^y(\mathbf{y}, \mathbf{z}) = \langle \phi_y(\mathbf{y}), \phi_y(\mathbf{z}) \rangle$  and denoting the resulting score vectors as  $\mathbf{k}^x(\mathbf{u}) = (k^x(\mathbf{x}_i, \mathbf{u}))_{i=1}^n$  and  $\mathbf{k}^y(\mathbf{v}) = (k^y(\mathbf{y}_i, \mathbf{v}))_{i=1}^n$  we can write (4) as

$$\rho_{\text{gradKCCA}} = \max_{\mathbf{u}, \mathbf{v}} \frac{\mathbf{k}^x(\mathbf{u})^T \mathbf{k}^y(\mathbf{v})}{\|\mathbf{k}^x(\mathbf{u})\|_2 \|\mathbf{k}^y(\mathbf{v})\|_2} \quad (5)$$

Note that the optimization problem (5) is the same as CCA (1) and KCCA (2) when  $\mathbf{k}^x$  and  $\mathbf{k}^y$  are the linear kernel functions. When  $\mathbf{k}^x$  and  $\mathbf{k}^y$  are non-linear kernels, the model constrains the projection directions to lie on the image of the feature map (see Section 4 for details). In terms of computation, (5) differs from KCCA by the time and space complexity of evaluating and maintaining the model, which is linear in the number of examples for gradKCCA, compared to the quadratic complexity of KCCA, while still allowing us to measure non-linear correlations in high-dimensional kernel-induced feature spaces.

We regularize the objective (5) by constraining the  $\ell_P$  norm (with  $P = 1, 2$ ) of  $\mathbf{u}$  and  $\mathbf{v}$ . Hence, with regularization, the final optimization problem is as follows

$$\begin{aligned} \rho_{\text{gradKCCAre}} = \max_{\mathbf{u}, \mathbf{v}} \rho(\mathbf{u}, \mathbf{v}) := & \frac{\mathbf{k}^x(\mathbf{u})^T \mathbf{k}^y(\mathbf{v})}{\|\mathbf{k}^x(\mathbf{u})\|_2 \|\mathbf{k}^y(\mathbf{v})\|_2} \\ \text{s.t.} \quad & \|\mathbf{u}\|_{P_u} \leq s_u \text{ and } \|\mathbf{v}\|_{P_v} \leq s_v \end{aligned} \quad (6)$$

where  $s_u$  and  $s_v$  are user-defined positive scalars. The solution to (6) results in sparse projection vectors by setting  $P_u = 1$  and/or  $P_v = 1$ . When the values of  $s_u$  and  $s_v$  are decreased the sparsity increases.

To maximize the correlation in (5), we use an alternating projected gradient approach (Algorithm 1). The algorithm computes the gradients  $\nabla \rho_{\mathbf{u}} = \frac{\partial \rho(\mathbf{u}, \mathbf{v})}{\partial \mathbf{u}}$  (resp.  $\nabla \rho_{\mathbf{v}}$ ) of the correlation  $\rho_{\text{gradKCCA}}$  with respect to  $\mathbf{u}$  and  $\mathbf{v}$  in alternating fashion. Each gradient computation entails computing the gradients of the kernel function  $\nabla_{\mathbf{u}} \mathbf{k}^x(\mathbf{u})$  between the training data and the current  $\mathbf{u}$ , (See Supplementary material), which can be done in time  $O((p+q)n)$  time per epoch,

where  $p$  and  $q$  are the numbers of variables in the two views, in the case of polynomial and RBF kernels. In addition, for large-scale data, a mini-batch approach can be used to facilitate optimization.

---

**Algorithm 1** gradKCCA
 

---

```

1: Input:  $\mathbf{X}, \mathbf{Y}, M$  (components),  $R$  (repetitions),
    $\delta$  (convergence limit),  $P_x$  and  $P_y$  (norms of  $\mathbf{u}$  and  $\mathbf{v}$ )
    $s_x$  and  $s_y$  ( $l_1$  or  $l_2$  norm constraints for  $\mathbf{u}$  and  $\mathbf{v}$ ),
    $d_x$  and  $d_y$  (hyperparameters for  $\mathbf{k}^x$  and  $\mathbf{k}^y$ )
2: Output:  $\mathbf{U}, \mathbf{V}$ 
3: for all  $m = \{1, 2, \dots, M\}$  do
4:   for all  $r = \{1, 2, \dots, R\}$  do
5:     Initialize  $\mathbf{u}_{mr}$  and  $\mathbf{v}_{mr}$ 
6:     Compute  $\mathbf{k}^x(\mathbf{u}), \mathbf{k}^y(\mathbf{v})$ 
7:     repeat
8:       Compute  $\rho_{\text{old}} = \rho(\mathbf{u}, \mathbf{v})$ 
9:       Compute  $\nabla \rho_{\mathbf{u}} = \frac{\partial \rho(\mathbf{u}, \mathbf{v})}{\partial \mathbf{u}}$ 
10:      Update  $\mathbf{u}_{mr} = \prod_{\|\cdot\|_{P_x} \leq s_x} (\mathbf{u}_{mr} + \gamma \nabla \rho_{\mathbf{u}})$ 
      (The step size  $\gamma$  determined by line search)
11:      Re-compute  $\mathbf{k}^x(\mathbf{u})$ 
12:      Compute  $\nabla \rho_{\mathbf{v}} = \frac{\partial \rho(\mathbf{u}, \mathbf{v})}{\partial \mathbf{v}}$ 
13:      Update  $\mathbf{v}_{mr} = \prod_{\|\cdot\|_{P_y} \leq s_y} (\mathbf{v}_{mr} + \gamma \nabla \rho_{\mathbf{v}})$ 
      (The step size  $\gamma$  determined by line search)
14:      Re-compute  $\mathbf{k}^y(\mathbf{v})$ 
15:      Compute  $\rho_{\text{current}} = \rho(\mathbf{u}, \mathbf{v})$ 
16:      until  $|\rho_{\text{old}} - \rho_{\text{current}}| / |\rho_{\text{old}} + \rho_{\text{current}}| < \delta$ 
17:       $\rho_r = \rho_{\text{current}}, \mathbf{u}_r = \mathbf{u}_{mr}, \mathbf{v}_r = \mathbf{v}_{mr}$ 
18:    end for
19:    Select  $r^* = \arg \max_r \rho_r$ 
20:    Store  $\mathbf{U}(:, m) = \mathbf{u}_{r^*}, \mathbf{V}(:, m) = \mathbf{v}_{r^*}$ 
21:    Deflate  $\mathbf{X}^{(m)}, \mathbf{Y}^{(m)}$  by  $\mathbf{U}(:, m)$  and  $\mathbf{V}(:, m)$ 
22:  end for
23: Return:  $\mathbf{U}, \mathbf{V}$ 
    
```

---

In Algorithm 1, the  $\ell_1$  or  $\ell_2$  norm constraints are implemented as follows. Let  $\prod_{\|\cdot\|_P \leq s}(\mathbf{u})$  denote the projection of  $\mathbf{u}$  onto an  $\ell_P$  norm ball. We solve

$$\begin{aligned} \prod_{\|\cdot\|_P \leq s}(\mathbf{u}) &= \min_{\mathbf{u}_{\text{project}}} \|\mathbf{u}_{\text{project}} - \mathbf{u}\|_2 \\ \text{s.t.} \quad &\|\mathbf{u}_{\text{project}}\|_P \leq s \end{aligned}$$

When  $P = 2$ ,  $\prod_{\|\cdot\|_P \leq s}(\mathbf{u}) = \frac{s\mathbf{u}}{\|\mathbf{u}\|_2}$ . When  $P = 1$ , we follow the method of  $\ell_1$  ball projection described in (Duchi et al., 2008).

The local convergence of Algorithm 1 follows from the theory presented in (Iusem, 2003). The step size of Algorithm 1 is found by Armijo line search, as described in Equations (13) and (14) in (Iusem, 2003). Following (Iusem (2003), Theorem 2), we claim that Algorithm 1 generates two sequences  $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k, \dots\}$  and  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k, \dots\}$  such that for all  $k$ ,  $\|\mathbf{u}_k\|_{P_x} \leq s_x$  and  $\|\mathbf{v}_k\|_{P_y} \leq s_y$ ,

all step sizes  $\gamma_k$  are well-defined. If  $\nabla \rho_{\mathbf{u}}(\mathbf{u}_k) \neq 0$  then  $\langle \nabla \rho_{\mathbf{u}}(\mathbf{u}_k), \text{Proj}(\mathbf{u}_k + \gamma_k \nabla \rho_{\mathbf{u}}(\mathbf{u}_k)) - \mathbf{u}_k \rangle > 0$  and if  $\nabla \rho_{\mathbf{v}}(\mathbf{v}_k) \neq 0$ , then  $\langle \nabla \rho_{\mathbf{v}}(\mathbf{v}_k), \text{Proj}(\mathbf{v}_k + \gamma_k \nabla \rho_{\mathbf{v}}(\mathbf{v}_k)) - \mathbf{v}_k \rangle > 0$ . Hence, at every step, the objective function is increasing if the gradient is nonzero. Assuming a locally concave function  $\rho(\mathbf{u}, \mathbf{v})$ , Algorithm 1 stops at some iteration  $k$ , in which case  $(\mathbf{u}, \mathbf{v})$  is a solution to Problem 6, or it generates infinite sequences  $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k, \dots\}$  and  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k, \dots\}$  that converge to a solution  $(\mathbf{u}^*, \mathbf{v}^*)$  of the problem.

### 3.1. Finding multiple canonical vectors

Given up to the  $m^{\text{th}}$  estimate of canonical projection  $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m\}$  and  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ , the following orthogonal canonical projection vectors can be obtained by solving (6) for  $\mathbf{u}_{m+1}$  and  $\mathbf{v}_{m+1}$  with additional constraint that  $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{m+1}\}$  and  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{m+1}\}$  are mutually orthogonal. In Algorithm 1, deflation (Mackey, 2009) is applied to obtain multiple multivariate relations. As an example for  $\mathbf{u}$ , upon finding the  $m^{\text{th}}$  estimate for  $\mathbf{u}_m$  from the current data matrix  $\mathbf{X}^{(m)}$ , the following orthogonal projection vector  $\mathbf{u}_{(m+1)}$  can be obtained by solving  $\mathbf{X}^{(m+1)} = \mathbf{X}^{(m)} - \frac{\mathbf{u}_m \mathbf{u}_m^T \mathbf{X}^{(m)T}}{\mathbf{u}_m^T \mathbf{u}_m}$ . The same procedure is applied on  $\mathbf{v}$ . In this way, the resulting sets of vectors  $\{\mathbf{u}_1, \mathbf{u}_2, \dots\}$  and  $\{\mathbf{v}_1, \mathbf{v}_2, \dots\}$  are mutually orthogonal. An alternative approach for finding orthogonal projection directions, is to look for uncorrelated projections (Hardoon et al., 2004). These can be obtained by solving (6) for  $\mathbf{u}_{m+1}$  and  $\mathbf{v}_{m+1}$  with additional constraint that the score vectors  $\{\mathbf{k}^x(\mathbf{u}_1), \dots, \mathbf{k}^x(\mathbf{u}_m)\}$  and  $\{\mathbf{k}^y(\mathbf{v}_1), \dots, \mathbf{k}^y(\mathbf{v}_m)\}$  are mutually orthogonal.

Instead of using the successive formulation of Algorithm 1, we can join the successive sub-problems into one and solve (7) for matrices  $\mathbf{U} \in \mathbb{R}^{p \times M}$  and  $\mathbf{V} \in \mathbb{R}^{q \times M}$  which contains  $M$  vectors  $(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_M)$  and  $(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_M)$ .

$$\begin{aligned} \max_{\mathbf{U}, \mathbf{V}} \quad &\sum_{i=1}^M \frac{\mathbf{k}^x(\mathbf{u}_i)^T \mathbf{k}^y(\mathbf{v}_i)}{\|\mathbf{k}^x(\mathbf{u}_i)\| \|\mathbf{k}^y(\mathbf{v}_i)\|} \\ \text{s.t.} \quad &\mathbf{k}^x(\mathbf{u}_i)^T \mathbf{k}^x(\mathbf{u}_j) = 0, \quad \forall i, j, i \neq j \\ &\mathbf{k}^y(\mathbf{v}_i)^T \mathbf{k}^y(\mathbf{v}_j) = 0, \quad \forall i, j, i \neq j \\ &\mathbf{k}^x(\mathbf{u}_i)^T \mathbf{k}^y(\mathbf{v}_j) = 0, \quad \forall i, j, i \neq j \end{aligned} \quad (7)$$

The non-linear equality constraints in (7) are considerably harder to solve than (6) with the deflation approach. One can soften constraints by adding them to the objective function using the concept of penalty method and then solve the unconstrained problem by gradient descent. We leave following this direction for future work.



#### 4. Relation between gradKCCA and KCCA

Let  $\alpha^*$  and  $\beta^*$  be the optimal solution of KCCA (2). Then the optimal canonical vectors (Kuss & Graepel, 2003) are

$$\mathbf{w}_x^* = \sum_i \alpha_i^* \phi_x(\mathbf{x}_i) \in \mathcal{L}(\phi_x, \mathbf{X}) \subseteq \mathcal{H}_x \text{ and}$$

$$\mathbf{w}_y^* = \sum_i \beta_i^* \phi_y(\mathbf{y}_i) \in \mathcal{L}(\phi_y, \mathbf{Y}) \subseteq \mathcal{H}_y.$$

where  $\mathcal{L}(\phi_x, \mathbf{X})$  and  $\mathcal{L}(\phi_y, \mathbf{Y})$  are spaces spanned by feature maps corresponding to training samples only.

Let  $\mathbf{u}^*$  and  $\mathbf{v}^*$  be the optimum solution for (5) and hence  $\phi_x(\mathbf{u}^*)$  and  $\phi_y(\mathbf{v}^*)$  are the optimal canonical vectors. Please note that,  $\phi_x(\mathbf{u}^*)$  (resp.  $\phi_y(\mathbf{v}^*)$ ) does not necessarily lie within the span  $\mathcal{L}(\phi_x, \mathbf{X})$  (resp.  $\mathcal{L}(\phi_y, \mathbf{Y})$ ) of the training data in view  $\mathbf{X}$  (resp. view  $\mathbf{Y}$ ).

First we show that KCCA (2) correlation  $\rho_{kcca}$  upper-bounds the correlation found by gradKCCA (5),  $\rho_{gradKCCA}$ .

**Lemma 4.1.**  $\rho_{gradKCCA} \leq \rho_{kcca}$

*Proof.* Considering the above notations, in general  $\phi_x(\mathbf{u}^*)$  or  $\phi_y(\mathbf{v}^*)$  may not lie in  $\mathcal{L}(\phi_x, \mathbf{X})$  and  $\mathcal{L}(\phi_y, \mathbf{Y})$ . If  $\phi_x(\mathbf{u}^*) \notin \mathcal{L}(\phi_x, \mathbf{X})$ ,  $\phi_x(\mathbf{u}^*)$  can be split into the following two components

$$\phi_x(\mathbf{u}^*) = \phi_x(\mathbf{u}^*)^{\parallel} + \phi_x(\mathbf{u}^*)^{\perp}$$

where  $\phi_x(\mathbf{u}^*)^{\parallel} \in \mathcal{L}(\phi_x, \mathbf{X})$  and  $\phi_x(\mathbf{u}^*)^{\perp}$  is orthogonal  $\mathcal{L}(\phi_x, \mathbf{X})$ , that is,  $\langle \phi_x(\mathbf{u}^*)^{\perp}, \phi(x_i) \rangle = 0$  for any  $i = 1, \dots, n$ . The score vector  $\mathbf{k}^x(\mathbf{u}^*)$  satisfies

$$\begin{aligned} \mathbf{k}^x(\mathbf{u}^*) &= \left[ \langle \phi_x(\mathbf{x}_i), \phi_x(\mathbf{u}^*) \rangle \right]_{i=1}^n = \\ &= \left[ \langle \phi_x(\mathbf{x}_i), \phi_x(\mathbf{u}^*)^{\parallel} \rangle \right]_{i=1}^n + \left[ \langle \phi_x(\mathbf{x}_i), \phi_x(\mathbf{u}^*)^{\perp} \rangle \right]_{i=1}^n \\ &= \left[ \langle \phi_x(\mathbf{x}_i), \phi_x(\mathbf{u}^*)^{\parallel} \rangle \right]_{i=1}^n + \mathbf{0}. \end{aligned}$$

In other words, the orthogonal part has no contribution to the score vector  $\mathbf{k}^x(\mathbf{u}^*)$ . An analogous fact holds true for the score vector  $\mathbf{k}^y(\mathbf{v}^*)$ . Consequently, the value of the canonical correlation  $\rho_{gradKCCA}$  is also not affected by the parts orthogonal to the span. Hence there exist some  $\alpha$  and  $\beta$  such that  $\phi_x(\mathbf{u}^*)^{\parallel} = \sum_i \alpha_i \phi_x(\mathbf{x}_i)$  and  $\phi_y(\mathbf{v}^*)^{\parallel} = \sum_i \beta_i \phi_y(\mathbf{y}_i)$  that achieve the same correlation as gradKCCA. Since KCCA (2) returns the optimal  $\alpha^*$  and  $\beta^*$  the claim follows.  $\square$

We note that above the dual variables for gradKCCA were used only as a proof technique, neither them nor the kernel matrix is needed in gradKCCA optimization.

Since for non-linear maps  $\phi_x$ ,  $\mathcal{L}(\phi_x, \mathbf{X})$  does not in general exactly coincide with the image of the feature map  $\phi_x$ , it

may be that some vectors in  $\mathbf{w}_x \in \mathcal{L}(\phi_x, \mathbf{X})$  and  $\mathbf{w}_y \in \mathcal{L}(\phi_y, \mathbf{Y})$ , the preimages  $\phi_x^{-1}(\mathbf{w}_x)$  and  $\phi_y^{-1}(\mathbf{w}_y)$  do not exist (Weston et al., 2004). In these cases, in general the (approximate) preimage  $\mathbf{u}$  for is chosen such that the squared distance of  $\mathbf{w}_x$  and  $\phi_x(\mathbf{u})$  is minimized.

$$\tilde{\mathbf{u}}(\mathbf{w}_x) = \arg \min_{\mathbf{u} \in \mathbb{R}^p} \|\mathbf{w}_x - \phi_x(\mathbf{u})\|_2 \quad (8)$$

$$\tilde{\mathbf{v}}(\mathbf{w}_y) = \arg \min_{\mathbf{v} \in \mathbb{R}^q} \|\mathbf{w}_y - \phi_y(\mathbf{v})\|_2 \quad (9)$$

The correlation achieved by computing the preimage and mapping back to the feature space is given by

$$\rho_{preimage} = \frac{\mathbf{k}^x(\tilde{\mathbf{u}}(\mathbf{w}_x))^{\top} \mathbf{k}^y(\tilde{\mathbf{v}}(\mathbf{w}_y))}{\|\mathbf{k}^x(\tilde{\mathbf{u}}(\mathbf{w}_x))\| \|\mathbf{k}^y(\tilde{\mathbf{v}}(\mathbf{w}_y))\|} \quad (10)$$

This preimage gives a lower-bound on  $\rho_{gradKCCA}$ :

**Lemma 4.2.**  $\rho_{gradKCCA} \geq \rho_{preimage}$

*Proof.* Follows from the fact that (5) gives the the maximum correlation for all  $\mathbf{u} \in \mathbb{R}^p$  and  $\mathbf{v} \in \mathbb{R}^q$ , including  $\mathbf{u} = \tilde{\mathbf{u}}(\mathbf{w}_x)$  and  $\mathbf{v} = \tilde{\mathbf{v}}(\mathbf{w}_y)$   $\square$

This preimage can be efficiently computed from an optimal KCCA solution. Namely, substituting the KCCA optimum  $\mathbf{w}_x^* = \sum_i \alpha_i^* \phi_x(x_i)$  into (8) we obtain

$$\begin{aligned} \tilde{\mathbf{u}}(\mathbf{w}_x^*) &= \arg \min_{\mathbf{u}} \left\| \sum_i \alpha_i^* \phi(\mathbf{x}_i) - \phi(\mathbf{u}) \right\|^2 = \\ &= \arg \min_{\mathbf{u}} \left( \left\| \sum_i \alpha_i^* \phi(\mathbf{x}_i) \right\|^2 - 2 \left\langle \sum_i \alpha_i^* \phi(\mathbf{x}_i), \phi(\mathbf{u}) \right\rangle \right. \\ &\quad \left. + \|\phi(\mathbf{u})\|^2 \right) \\ &= \arg \min_{\mathbf{u}} \alpha^{*T} \mathbf{K}^x \alpha^* - 2 \alpha^{*T} \mathbf{k}^x(\mathbf{u}) + \mathbf{k}^x(\mathbf{u}, \mathbf{u}) \end{aligned}$$

which can be solved by e.g. by gradient descent on  $\mathbf{u}$ . We use this approach in our experiment section to compute preimages for KCCA, and name it as  $KCCA_{preimage}$ .

Using the above preimage problem, we can further bound the distance of gradKCCA from the KCCA optimum. Let us define the bounds on preimage errors as

$$\begin{aligned} B_{\mathbf{X}} &= \max_{\mathbf{w}_x \in \mathcal{L}(\phi_x, \mathbf{X})} \|\mathbf{w}_x - \phi_x(\tilde{\mathbf{u}}(\mathbf{w}_x))\|_2 \\ B_{\mathbf{Y}} &= \max_{\mathbf{w}_y \in \mathcal{L}(\phi_y, \mathbf{Y})} \|\mathbf{w}_y - \phi_y(\tilde{\mathbf{v}}(\mathbf{w}_y))\|_2. \end{aligned} \quad (11)$$

**Theorem 4.3.** Let us further assume that norm in  $\mathcal{H}_x$  and  $\mathcal{H}_y$  are upper bounded by  $M_x$  and  $M_y$ , that is,

$$\begin{aligned} \forall \phi_x(\mathbf{u}) \in \mathcal{H}_x, \|\phi_x(\mathbf{u})\| &\leq M_x \\ \text{and } \forall \phi_y(\mathbf{v}) \in \mathcal{H}_y, \|\phi_y(\mathbf{v})\| &\leq M_y. \end{aligned} \quad (12)$$

Then,

$$\rho_{gradKCCA} \geq \rho_{preimage} \geq \rho_{kcca} - \left( \frac{B_y}{M_y} + \frac{B_x}{M_x} \right)$$

*Proof.* See Section 2 in the Supplementary Material.  $\square$

## 5. Experiments

### 5.1. Predictive performance

We assess the performance of gradKCCA, and compare it with KCCA, DCCA, RCCA, KNOI and SCCA-HSIC, in terms of robustness to noise and scalability. In the experiments, we test two different sets of input relations, monotone algebraic relations and non-monotone trigonometric relations. For the monotone algebraic relations, we use the linear kernel for gradKCCA and KCCA. For the non-monotone trigonometric relations we use the quadratic kernel for both. The code is available on <https://github.com/aalto-ics-kepaco/gradKCCA>.

For the other kernel-based methods, that is RCCA, KNOI, and SCCA-HSIC, we only apply the Gaussian kernel, which is a requirement for SCCA-HSIC due to its universality (Gretton et al., 2008), and which has been applied previously together with random Fourier features for RCCA and KNOI (Lopez-Paz et al., 2014; Wang & Livescu, 2016).

**Experiments on monotone relations.** In the experiments, where the linear kernels are tested for gradKCCA and KCCA, we apply the following monotone algebraic relations. We generate polynomial relations of the form  $f(x) = x^d$ , where  $x \in \mathbb{R}_+$  and  $d = 1, 2, 3$ , that is linear, quadratic, and cubic relations, respectively. We also generate a transcendental relation, the exponential relation, of the form  $f(x) = \exp(x)$ , where  $x \in \mathbb{R}_+$ . We consider a two-to-two relation, that is, the two first columns of  $\mathbf{Y}$  and  $\mathbf{X}$  satisfy  $\mathbf{y}^1 + \mathbf{y}^2 = f(\mathbf{x}^1 + \mathbf{x}^2) + \boldsymbol{\xi}$  where  $\boldsymbol{\xi} \sim \mathcal{N}(0, 0.05)$  denotes a vector of normal distributed noise. The variables in the columns of  $\mathbf{X}$  and  $\mathbf{Y}$  are generated from a random uni-variate uniform distribution,  $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^p \sim U[0, 1]$  and  $\mathbf{y}^3, \mathbf{y}^4, \dots, \mathbf{y}^q \sim U[0, 1]$ .

**Experiments on non-monotone relations.** As non-monotone relations, we generate the following trigonometric relations that cannot be approximated using a linear kernel. Let every entry of  $\boldsymbol{\theta} \in \mathbb{R}^n \sim U[-2\pi, 2\pi]$ . The variables, in the columns of  $\mathbf{X}$  and  $\mathbf{Y}$ , are generated by

$$\begin{aligned} \mathbf{X} &= \left( \sin \frac{\mathbf{x}_2}{a} + \boldsymbol{\xi} \quad \mathbf{x}_2 \quad \cdots \quad \boldsymbol{\theta}_p \right) \\ \mathbf{Y} &= \left( \cos \frac{\mathbf{x}_2}{b} + \boldsymbol{\xi} \quad \cos \frac{\mathbf{x}_2}{c} + \boldsymbol{\xi} \quad \cdots \quad \boldsymbol{\theta}_q \right) \end{aligned}$$

where  $\mathbf{x}_2$  denotes the second variable in view  $\mathbf{X}$  and  $\boldsymbol{\xi} \sim \mathcal{N}(0, 0.05)$  denotes a vector of normal distributed noise. In the  $\mathbf{X}$  view, the first and second variables are related, and the first and second variables in the  $\mathbf{Y}$  view are related with the second variable of  $\mathbf{X}$  view. We generate three different variants of this setting. In the first relation, we put  $a = -2$ ,  $b = 3$ , and  $c = 4$ . In the second relation, we put  $a = -2$ ,  $b = -4$ , and  $c = -6$ . In the third relation, we put  $a = -0.5$ ,  $b = 2$ , and  $c = -2$ .

**Hyperparameter settings.** In all of the experiments, we randomly select 500 examples from the simulated training dataset for hyperparameter tuning by repeated nested cross-validation. For gradKCCA and SCCA-HSIC, we tune the  $\ell_1$  sparsity constraints. For KCCA, we tune the regularization hyperparameters. For DCCA, we tune the neural network structure, that is the number of layers and the number of hidden units. For RCCA and KNOI, we tune the number of random Fourier features.

**Evaluation metrics.** We assess the performances of the methods in terms of test correlation, F1 score, and Area Under the Curve (AUC). The test correlations are computed using (5), where the learned  $\mathbf{u}$  and  $\mathbf{v}$  are applied on test examples.

For gradKCCA, SCCA-HSIC and  $\text{KCCA}_{\text{preimage}}$ , we assess the correctness of  $\mathbf{u}$  and  $\mathbf{v}$  using the F1 score and AUC. For  $\text{KCCA}_{\text{preimage}}$ , we use the output of the preimage algorithm as the predicted  $\mathbf{u}$  and  $\mathbf{v}$ . F1 score is defined as  $\text{F1} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$ , where TP, FP, and FN denote the true positives, false positives, and false negatives respectively. In our case, the labels refer to the zero/non-zero coefficient values in the estimated  $\mathbf{u}$  and  $\mathbf{v}$  and the respective ground truth vectors. For the computations of the two metrics, the estimated  $\mathbf{u}$  and  $\mathbf{v}$  are binarized such that all entries less than or equal to 0.05 are put to zero and all the rest to 1.

#### 5.1.1. VARYING THE GROUND TRUTH SPARSITY

We evaluate the robustness of gradKCCA,  $\text{KCCA}_{\text{preimage}}$ , KCCA, DCCA, KNOI and SCCA-HSIC when the sparsity of the ground truth, that is the number of noise variables, increases. We do not include RCCA in these small-scale experiments since it results in the same kernel canonical correlation as KCCA at a sufficient number of random Fourier features.

For every multivariate function, given in Section 5.1, we repeat the following scheme 10 times. We generate two data matrices  $\mathbf{X}$  and  $\mathbf{Y}$  of sizes  $n \times p$  and  $n \times q$ , where  $n = 1000$ . The dimensions tested are  $p = q = \{5, 10, 50, 100\}$ .

**Monotone relations.** The results are shown in Figure 1. Since a linear kernel function is applied, gradKCCA and KCCA perform similarly, in terms of training and test correlation. The F1 score and AUC show that gradKCCA is slightly more robust to noise than KCCA. This could be due to the constraint that the projection direction lies on the image of the feature map. DCCA, KNOI, and SCCA-HSIC do not tolerate noise as well as gradKCCA and KCCA. For KNOI, this could be explained by the stochastic approximation algorithm that uses mini-batches of the training examples.

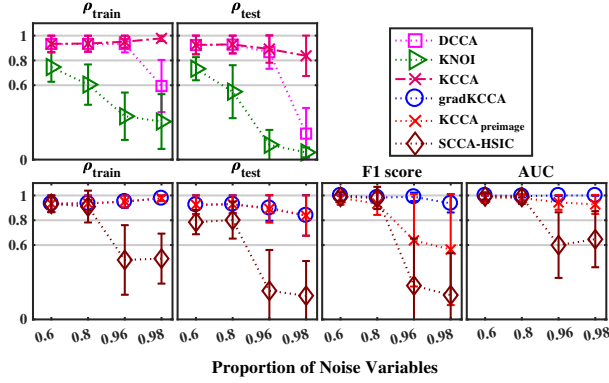


Figure 1. Experiment 5.1.1 on monotone relations. DCCA, KNOI, and KCCA and gradKCCA, KCCA<sub>preimage</sub>, and SCCA-HSIC are shown on upper and lower rows respectively.

**Non-monotone relations.** The results are shown in Figure 2. In terms of test correlations, the quadratic kernel gradKCCA and SCCA-HSIC perform equally well. However, gradKCCA is more accurate than SCCA-HSIC in identifying the related variables in terms of F1 score. As has been shown in Section 4, in the case of higher degree kernels, the canonical correlation using the computed the preimages of KCCA is lower than the kernel canonical correlation. KCCA, DCCA and KNOI are shown not to be tolerant to irrelevant noise variables in this context.

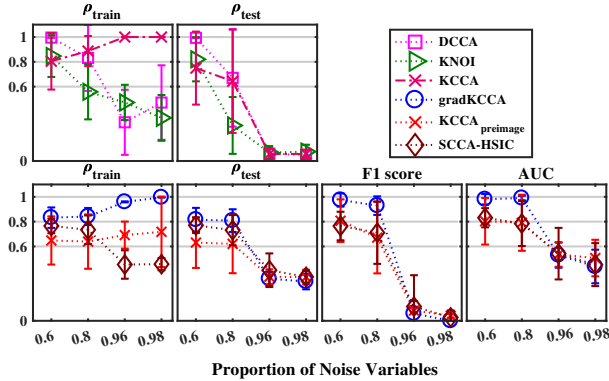


Figure 2. Experiment 5.1.1 on non-monotone relations. See the caption of Figure 1.

### 5.1.2. VARYING THE SAMPLE SIZE

In the scalability experiments, we compare gradKCCA with the large-scale CCA variants, that is DCCA, RCCA, KNOI, and Nyström approximated SCCA-HSIC. For every multivariate function, given in Section 5.1, we repeat the following scheme five times. We generate two data matrices  $X$  and

$Y$  of sizes  $n \times p$  and  $n \times q$ , where  $n = \{10^3, 10^4, 10^5, 10^6\}$ . The number of inducing variables for the Nyström approximated SCCA-HSIC is set to  $\{0.7, 0.1, 0.05\}$  for the first three sample sizes respectively. The dimensions are  $p = q = 10$ .

**Monotone relations.** The results are shown in Figure 3. In terms of training and test (kernel) canonical correlations, gradKCCA, RCCA, and DCCA perform equally well. However, in terms of running time gradKCCA outperforms all the methods. SCCA-HSIC is the slowest method, and is not run for one million examples.

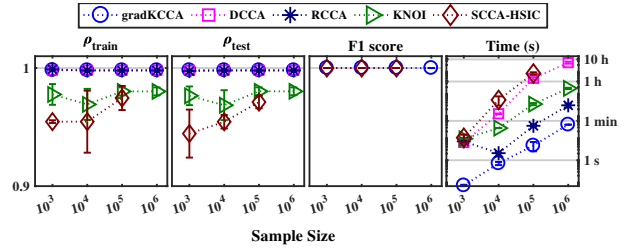


Figure 3. Experiment 5.1.2 on monotone relations. The performances of gradKCCA, DCCA, RCCA, KNOI, and SCCA-HSIC are shown when the sample sizes  $n = \{10^3, 10^4, 10^5, 10^6\}$  are tested. Due to the slow computation time, SCCA-HSIC is only run until  $10^5$ .

**Non-monotone relations.** The results are shown in Figure 4. As in the experiment on monotone relations, gradKCCA is the fastest method. DCCA and KNOI extract non-monotone relations better at larger sample sizes. In contrast to the previous experiment, gradKCCA is better in extracting the relations than SCCA-HSIC, which is seen both in terms of test correlation and F1 score.

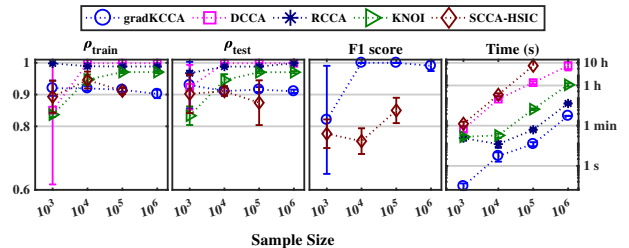


Figure 4. Experiment 5.1.2 on non-monotone relations. See the caption of Figure 3.

## 5.2. Real-world datasets

We evaluate the performances of gradKCCA, DCCA, RCCA, KNOI, and SCCA-HSIC in terms of the training

Table 1. MNIST dataset.

	$\rho_{\text{TRAIN}}$	$\rho_{\text{TEST}}$	TIME (S)
GRADKCCA	$0.955 \pm 0.001$	$0.952 \pm 0.001$	$56 \pm 6$
DCCA	$0.941 \pm 0.010$	$0.943 \pm 0.003$	$4578 \pm 203$
RCCA	$0.949 \pm 0.001$	$0.949 \pm 0.010$	$78 \pm 13$
KNOI	$0.950 \pm 0.001$	$0.950 \pm 0.005$	$878 \pm 62$
SCCA-HSIC	$0.912 \pm 0.020$	$0.934 \pm 0.006$	$5611 \pm 193$

and test correlation of the first CCA component and computation time on the publicly available MNIST handwritten digits (LeCun, 1998) and MediaMill (Snoek et al., 2006) datasets. We extract the first component 10 times to obtain an average. We tune the hyperparameters in similar manner as in Section 5.1.

**MNIST handwritten digits.** The MNIST handwritten digits dataset contains 60 000 training and 10 000 testing images of handwritten digits. As in (Andrew et al., 2013), we analyze the relations of the left and right halves of images of the handwritten digits. Every image consists of  $28 \times 28$  matrix of pixels. Every pixel represents a grey scale value ranging from one to 256. The 14 left and right columns are separated to form the two views, making 392 features in each view.

The training and test correlations of the views, and computation times, are given in Table 1. gradKCCA outperforms in terms of speed, and finds a test correlation of highest value. KNOI is the second best in terms of test correlation, and RCCA the third best, with the second fastest computation time. The complexity of RCCA is  $O((m_x^2 + m_y^2)n)$  where  $m_x$  and  $m_y$  are dimensions of the randomized feature map and the complexity of calculating feature maps are  $O(pm_xn)$  and  $O(qm_yn)$ . This explains why RCCA is slower than gradKCCA in practice.

**MediaMill dataset.** The MediaMill dataset consists of 43907 examples that correspond to keyframes taken from video shots. Every keyframe is represented by 120 visual features and 101 text annotations, or labels, respectively.

The training and test correlations between the visual and annotation views, together with the computation times, are given in Table 2. As in the MNIST experiment, gradKCCA is the fastest. RCCA is the second best in terms of speed and KNOI is the second best in terms of test correlation.

## 6. Discussion

Our experiments showed that the gradKCCA model is orders of magnitude faster to optimize than competing non-linear CCA models. On real-world datasets, gradKCCA obtained top predictive performance, and again proved to

Table 2. MediaMill Dataset.

	$\rho_{\text{TRAIN}}$	$\rho_{\text{TEST}}$	TIME (S)
GRADKCCA	$0.666 \pm 0.004$	$0.657 \pm 0.007$	$8 \pm 4$
DCCA	$0.643 \pm 0.005$	$0.633 \pm 0.003$	$1280 \pm 112$
RCCA	$0.633 \pm 0.001$	$0.626 \pm 0.005$	$23 \pm 9$
KNOI	$0.652 \pm 0.001$	$0.645 \pm 0.003$	$218 \pm 73$
SCCA-HSIC	$0.627 \pm 0.004$	$0.625 \pm 0.002$	$1804 \pm 143$

be the fastest method by a large margin. In our experiments with artificial data, we also showed that gradKCCA is very robust to a large amount respect to irrelevant, noise variables in the data. Also, the method obtains good performance already with modest amount of data, compared to several of the competing methods.

The gradKCCA model can be seen as a variant of KCCA that gives access to the preimages of the canonical projection directions. Our results in Section 4 show that the correlations on the training set are upper bounded by KCCA and lower-bounded by correlations that are obtained after taking the preimage of the KCCA solution. In our experiments, these relations were generally seen to translate to the test set as well. Thus, gradKCCA can be seen to give a good compromise in situations where the model’s interpretation in terms of relevant variables is of interest.

The general approach used here that removes the need to maintain a kernel matrix by moving the optimization to the preimage space is not limited to CCA models. We envision that many other kernel methods could be scaled up to big data using the same principle. On the other hand, it is clear that the new model puts more burden for finding a good kernel so that the requirement of the kernel space model having an exact preimage does not become too restrictive.

## 7. Conclusions

In this paper, we have put forward a new method for finding canonical correlations in paired datasets. The proposed method, gradKCCA, adapts KCCA by requiring the projection directions in the feature space to have preimages in the original space. This modeling choice has the by-product that a kernel matrix does not need to be computed or maintained. This allows fast optimization through gradient approaches, which scales up to big datasets effortlessly. Moreover due the approach, the sparsity projection directions can be controlled which gives the method robustness against irrelevant variables.

## Acknowledgements

This work has been in part supported by Academy of Finland grants 310107 (MACOME) and 313268 (TensorBiomed).



## References

- Andrew, G., Arora, R., Bilmes, J., and Livescu, K. Deep canonical correlation analysis. In *ICML*, pp. 1247–1255, 2013.
- Bach, F. and Jordan, M. Kernel independent component analysis. *JMLR*, 3(Jul):1–48, 2002.
- Chang, B., Kruger, U., Kustra, R., and Zhang, J. Canonical correlation analysis based on hilbert-schmidt independence criterion and centered kernel target alignment. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pp. 316–324, 2013.
- Chu, D., Liao, L.-Z., Ng, M. K., and Zhang, X. Sparse kernel canonical correlation analysis. In *Proceedings of the International MultiConference of Engineers and Computer Scientists*, volume 1, 2013.
- Duchi, J., Shalev-Shwartz, S., Singer, Y., and Chandra, T. Efficient projections onto the  $l_1$ -ball for learning in high dimensions. In *Proceedings of the 25th international conference on Machine learning*, pp. 272–279. ACM, 2008.
- Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. Measuring statistical dependence with hilbert-schmidt norms. In *ALT*, volume 16, pp. 63–78. Springer, 2005.
- Gretton, A., Fukumizu, K., Teo, C., Song, L., Schölkopf, B., and Smola, A. A kernel statistical test of independence. In *Advances in neural information processing systems*, pp. 585–592, 2008.
- Hardoon, D., Szedmak, S., and Shawe-Taylor, J. Canonical correlation analysis: An overview with application to learning methods. *Neural computation*, 16(12):2639–2664, 2004.
- Hotelling, H. The most predictable criterion. *Journal of educational Psychology*, 26(2):139, 1935.
- Hotelling, H. Relations between two sets of variates. *Biometrika*, 28(3):321–377, 1936.
- Iusem, A. On the convergence properties of the projected gradient method for convex optimization. *Computational & Applied Mathematics*, 22(1):37–52, 2003.
- Kuss, M. and Graepel, T. The geometry of kernel canonical correlation analysis. *Technical Report 108, Max-Planck Institute for Biol. Cybernetics*, 2003.
- LeCun, Y. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Lopez-Paz, D., Sra, S., Smola, A., Ghahramani, Z., and Schölkopf, B. Randomized nonlinear component analysis. In *International Conference on Machine Learning*, pp. 1359–1367, 2014.
- Ma, Z., Lu, Y., and Foster, D. Finding linear structure in large datasets with scalable canonical correlation analysis. In *International Conference on Machine Learning*, pp. 169–178, 2015.
- Mackey, L. W. Deflation methods for sparse pca. In *Advances in neural information processing systems*, pp. 1017–1024, 2009.
- Rahimi, A. and Recht, B. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pp. 1177–1184, 2008.
- Snoek, C. G., Worring, M., Van Gemert, J. C., Geusebroek, J.-M., and Smeulders, A. W. The challenge problem for automated detection of 101 semantic concepts in multimedia. In *Proceedings of the 14th ACM international conference on Multimedia*, pp. 421–430. ACM, 2006.
- Urtio, V., Bhadra, S., and Rousu, J. Sparse non-linear cca through hilbert-schmidt independence criterion. In *2018 IEEE International Conference on Data Mining (ICDM)*, pp. 1278–1283, Nov 2018a. doi: 10.1109/ICDM.2018.00172.
- Urtio, V., Monteiro, J., Kandola, J., Shawe-Taylor, J., Fernandez-Reyes, D., and Rousu, J. A tutorial on canonical correlation methods. *ACM Computing Surveys (CSUR)*, 50(6):95:1–95:33, Jan 2018b. ISSN 0360-0300.
- Wang, W. and Livescu, K. Large-scale approximate kernel canonical correlation analysis. In *Proceedings of the 4th International Conference on Learning Representations (ICLR)*, 2016.
- Wang, W., Arora, R., Livescu, K., and Srebro, N. Stochastic optimization for deep cca via nonlinear orthogonal iterations. In *Communication, Control, and Computing (Allerton), 2015 53rd Annual Allerton Conference on*, pp. 688–695. IEEE, 2015.
- Weston, J., Schölkopf, B., and Bakir, G. H. Learning to find pre-images. In *Advances in neural information processing systems*, pp. 449–456, 2004.
- Williams, C. and Seeger, M. Using the nyström method to speed up kernel machines. In *Adv Neural Inf Process Syst*, number EPFL-CONF-161322, pp. 682–688, 2001.
- Zhang, Q., Filippi, S., Gretton, A., and Sejdinovic, D. Large-scale kernel methods for independence testing. *Stat Comput*, 28(1):113–130, 2018.