

Kernel analysis on Grassmann manifolds for action recognition

Mehrtash T. Harandi^{a,*}, Conrad Sanderson^b, Sareh Shirazi^c, Brian C. Lovell^d

^aNICTA, Locked Bag 8001, Canberra ACT 2601, Australia and The University of Queensland, School of ITEE, QLD 4072, Australia

^bNICTA, PO Box 6020, St Lucia QLD 4067, Australia and Queensland University of Technology (QUT), Brisbane QLD 4000, Australia

^cNICTA, PO Box 6020, St Lucia QLD 4067, Australia and The University of Queensland, School of ITEE, QLD 4072, Australia

^dThe University of Queensland, School of ITEE, QLD 4072, Australia

A B S T R A C T

Modelling video sequences by subspaces has recently shown promise for recognising human actions. Subspaces are able to accommodate the effects of various image variations and can capture the dynamic properties of actions. Subspaces form a non-Euclidean and curved Riemannian manifold known as a Grassmann manifold. Inference on manifold spaces usually is achieved by embedding the manifolds in higher dimensional Euclidean spaces. In this paper, we instead propose to embed the Grassmann manifolds into reproducing kernel Hilbert spaces and then tackle the problem of discriminant analysis on such manifolds. To achieve efficient machinery, we propose graph-based local discriminant analysis that utilises within-class and between-class similarity graphs to characterise intra-class compactness and inter-class separability, respectively. Experiments on KTH, UCF Sports, and Ballet datasets show that the proposed approach obtains marked improvements in discrimination accuracy in comparison to several state-of-the-art methods, such as the kernel version of affine hull image-set distance, tensor canonical correlation analysis, spatial-temporal words and hierarchy of discriminative space-time neighbourhood features.

1. Introduction

The goal of human action recognition is to automatically analyse and recognise what action is being undertaken in a given video, with one or more persons performing an action. Applications include content-based video analysis, security and surveillance, human-computer interaction, and animation (Turaga et al., 2008; Weinland et al., 2011). Subspace-based approaches, which are able to accommodate the effects of a wide range of image variations, have recently shown promising results for action recognition (Kim and Cipolla, 2009; Turaga and Chellappa, 2009; O'Hara et al., 2012). Moreover, subspaces can also capture the dynamic properties of videos (Turaga et al., 2011).

Subspaces form non-Euclidean and curved Riemannian manifolds known as Grassmann manifolds, allowing a video or an image-set to be conveniently represented as a point on a Grassmann manifold. Recent studies show that better performance can be achieved when the geometry of Riemannian spaces is explicitly considered (Hamm and Lee, 2008; Tuzel et al., 2008; Subbarao and Meer, 2009; Lui, 2012; Harandi et al., 2011; Turaga et al., 2011).

Inference on manifold spaces can be achieved by embedding the manifolds in higher dimensional Euclidean spaces, which can be

considered as flattening the manifolds. A popular choice for embedding manifolds is through considering tangent spaces (Tuzel et al., 2008; Turaga et al., 2011). Two notable examples are the pedestrian detection system by Tuzel et al. (2008) and non-linear mean shift (Comaniciu and Meer, 2002) by Subbarao and Meer (2009). Nevertheless, flattening manifolds through tangent spaces is not without drawbacks. For example, the distance on a tangent space between two arbitrary points is generally not equal to the true geodesic distance,¹ which may lead to inaccurate modelling (Harandi et al., 2012). An alternative school of thought omits the use of tangent spaces and instead embeds Grassmann manifolds into Reproducing Kernel Hilbert Spaces (RKHS) (Shawe-Taylor and Cristianini, 2004) through dedicated Grassmann kernels (Hamm and Lee, 2008; Harandi et al., 2011; Shirazi et al., 2012). This in turn opens the door for employing many kernel-based machine learning algorithms (Shawe-Taylor and Cristianini, 2004).

Inference via discriminant analysis (DA) on Grassmann manifolds has been recently explored (Hamm and Lee, 2008; Wang and Shi, 2009). Given subspaces that are represented as points on a Grassmann manifold, Grassmann discriminant analysis (GDA) maps them to RKHS, such that a measure of discriminatory power in the induced RKHS is maximised. While GDA has shown promising results in (Hamm and Lee, 2008; Wang and Shi, 2009), the

* Corresponding author at: NICTA, Locked Bag 8001, Canberra ACT 2601, Australia and The University of Queensland, School of ITEE, QLD 4072, Australia. Tel.: +61 733008673.

E-mail address: mehrtash.harandi@nicta.com.au (M.T. Harandi).

¹ The geodesic distance takes into account the curvature of manifolds; an example is the distance between two points on a sphere.



Fig. 1. (a) Examples of a hand-waving action. (b) Basis vectors for a subspace of order three, modelling the entire action; the subspace is a point on a Grassmann manifold.

conventional formalism of DA suffers from not being able to take into account the local structure of data (Chen et al., 2007). For example, multi-modal classes and outliers can adversely affect the discrimination and/or generalisation ability of models based on conventional DA.

Contributions. In this work² we first present two methods of representing human actions on Grassmann manifolds. For the purposes of action recognition, we then extend our preliminary study on an enhanced form of GDA (Harandi et al., 2011), based on Grassmann kernels and a graph-embedding framework (Yan et al., 2007). We also show that conventional GDA (Hamm and Lee, 2008) can be seen as a special case of the enhanced graph-embedding based approach. Thorough experiments on the KTH (Schuldt et al., 2004), UCF Sports (Rodriguez et al., 2008) and Ballet (Wang and Mori, 2009) datasets, which include various realistic challenges such as background clutter, partial occlusion, changes in viewpoint, scale and illumination, and complexity of motion, show that the proposed Grassmann graph-embedding discriminant analysis (GGDA) approach obtains notable improvements in discrimination accuracy in comparison to several state-of-the-art methods. This includes the original GDA (Hamm and Lee, 2008), kernel version of affine hull image-set distance (Cevikalp and Triggs, 2010), tensor canonical correlation analysis (Kim and Cipolla, 2009), spatial-temporal words (Niebles et al., 2008) and hierarchy of discriminative space-time neighbourhood features (Kovashka and Grauman, 2010).

We continue the paper as follows. Section 2 presents various ways of representing action videos by linear subspaces. Section 3 reviews Grassmann geometry, which leads to Section 4 which presents the Grassmann graph-embedding discriminant analysis approach. In Section 5 we compare the performance of the proposed method with previous approaches on several datasets. The main findings and possible future directions are summarised in Section 6.

2. Modelling actions by linear subspaces

Let us define a video as an ordered collection of images with time-stamps (temporal information), and an image-set as an orderless collection of images. Actions can be modelled as linear subspaces through image-sets, or through linear dynamic systems that take into account the temporal information. We overview both methodologies in the following subsections.

2.1. Modelling of image-sets

An image-set $\mathbb{F} = \{\mathbf{f}_i\}_{i=1}^N; \mathbf{f}_i \in \mathbb{R}^n$, where \mathbf{f}_i is the vectorised representation of frame i , can be represented as a subspace (and hence as a point on a Grassmann manifold) through any orthogonalisation procedure like Singular Value Decomposition (SVD). More specifically, let $\mathbb{F} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ be the SVD of \mathbb{F} . The first p columns of \mathbf{U} represent an optimised subspace of order p (in the mean square sense) for \mathbb{F} and can be seen as a point on manifold $\mathcal{G}n, p$. Intuitively, modelling an action by a subspace as described here

can be understood as a low dimensional and compact representation by a set of basis vectors, in which the appearance of action could be effectively reconstructed by linearly combining the basis vectors. See Fig. 1 for an example of modelling a hand-waving action sequence by a subspace of order three.

Modelling image-sets by linear subspaces has been shown to deliver improved performance in the presence of practical issues such as misalignment as well as variations in pose and illumination (Wolf and Shashua, 2003; Hamm and Lee, 2008; Harandi et al., 2011). Modelling of actions by image-sets can be sufficient provided that the order in which the action is performed is not very relevant to decision making. While this assumption sounds restrictive, in many practical situations this might indeed be the case. As an example, it is possible to differentiate riding a horse from jogging without having temporal information. Nevertheless, a recent study (Li et al., 2011) shows that an extended type of image-set, obtained through a block Hankel matrix formalism, can capture the temporal information.

2.2. Modelling of linear dynamic systems

A video can be represented by an Auto Regressive and Moving Average (ARMA) model to explicitly take into account dynamics and temporal information. More specifically, a set of ordered images $\{\mathbf{f}(t)\}_{t=1}^{\tau}; \mathbf{f}(t) \in \mathbb{R}^n$ can be seen as the output of an ARMA process as:

$$\mathbf{f}(t) = \mathbf{C}\mathbf{z}(t) + \mathbf{w}(t), \quad \mathbf{w}(t) \sim \mathcal{N}(0, \mathbf{R}) \quad (1)$$

$$\mathbf{z}(t+1) = \mathbf{A}\mathbf{z}(t) + \mathbf{v}(t), \quad \mathbf{v}(t) \sim \mathcal{N}(0, \mathbf{Q}) \quad (2)$$

In (1) and (2), $\mathbf{z}(t) \in \mathbb{R}^p$ is the latent state vector at time t , $\mathbf{A} \in \mathbb{R}^{p \times p}$ and $\mathbf{C} \in \mathbb{R}^{n \times p}$ are the transition and measurement matrices, respectively, while \mathbf{w} and \mathbf{v} are noise components modelled as normal distributions with zero mean and covariance matrices $\mathbf{R} \in \mathbb{R}^{n \times n}$ and $\mathbf{Q} \in \mathbb{R}^{p \times p}$, respectively. The order of the system is given by p , while n is the number of features in a frame of the sequence. Loosely speaking, one advantage of the ARMA model is that it decouples the appearance of the spatio-temporal data (modelled by \mathbf{C}) from the dynamics (represented by \mathbf{A}).

To estimate the transition and measurement matrices (Turaga et al., 2011), we define $\mathbf{F}_{\tau} = [\mathbf{f}(1)|\mathbf{f}(2)|\dots|\mathbf{f}(\tau)]$, where the symbol $|$ denotes horizontal concatenation of vectors, as the feature matrix for time indices $1, 2, \dots, \tau$. The estimated transition ($\hat{\mathbf{A}}$) and measurement ($\hat{\mathbf{C}}$) matrices can then be obtained via the SVD of $\mathbf{F}_{\tau} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, as follows:

$$\hat{\mathbf{A}} = \mathbf{\Sigma}\mathbf{V}^T\mathbf{D}_1\mathbf{V}(\mathbf{V}^T\mathbf{D}_2\mathbf{V})^{-1}\mathbf{\Sigma}^{-1} \quad (3)$$

$$\hat{\mathbf{C}} = \mathbf{U} \quad (4)$$

where

$$\mathbf{D}_1 = \begin{bmatrix} \mathbf{0}_{\tau-1}^T & 0 \\ \mathbf{I}_{(\tau-1) \times (\tau-1)} & \mathbf{0}_{\tau-1} \end{bmatrix} \quad \text{and} \quad \mathbf{D}_2 = \begin{bmatrix} \mathbf{I}_{(\tau-1) \times (\tau-1)} & \mathbf{0}_{\tau-1} \\ \mathbf{0}_{\tau-1}^T & 0 \end{bmatrix}$$

ARMA models can be compared based on the subspace angles between the column-spaces of their observability matrices (Cock and Moor, 2002). The extended observability matrix of an ARMA model is given by $\mathbf{O}_{\infty} = [\mathbf{C}^T|(\mathbf{C}\mathbf{A})^T|(\mathbf{C}\mathbf{A}^2)^T|\dots|(\mathbf{C}\mathbf{A}^{n-1})^T|\dots]^T$ and is

² This work is somewhat related to (Shirazi et al., 2012), where the problem of clustering on Grassmann manifolds is explored. In the method presented here, clustering is not performed.

usually approximated by the finite observability matrix $\mathbf{O}_p = [\mathbf{C}^T | (\mathbf{C}\mathbf{A})^T | (\mathbf{C}\mathbf{A}^2)^T | \dots | (\mathbf{C}\mathbf{A}^{p-1})^T]^T$ (Turaga et al., 2011).

To represent a video on a Grassmann manifold, the finite observability matrix of the ARMA model is estimated as described above. The subspace spanned by the columns of \mathbf{O}_p (obtained by SVD or any other orthogonalisation procedure) corresponds to a point on the Grassmann manifold $\mathcal{G}_{n,p}$.

3. Grassmann geometry

Without delving too deeply into differential geometry and related topics, a Riemannian manifold \mathcal{M} is a differentiable and smooth manifold, endowed with a Riemannian metric that allows us to extend the notion of lengths and angles from familiar Euclidean spaces to the curved and non-flat space of \mathcal{M} . The geodesic distance between two points $\mathbf{X}, \mathbf{Y} \in \mathcal{M}$, denoted by $d_g(\mathbf{X}, \mathbf{Y})$, is defined as the minimum length over all possible smooth curves between \mathbf{X} and \mathbf{Y} . A geodesic curve is a curve that locally minimises the distance between points.

Subspaces form a special class of Riemannian manifolds known as Grassmann manifolds. Formally, Grassmann manifolds are defined as quotient spaces of orthogonal group³ $\mathbb{O}(n)$. A quotient space⁴ of a manifold, intuitively speaking, is the result of ‘‘gluing together’’ certain points of the manifold.

Definition 1. Grassmann manifold $\mathcal{G}_{n,p}$ is a quotient space of the orthogonal group $\mathbb{O}(n)$ and is defined as the set of p -dimensional linear subspaces of \mathbb{R}^n . Points on the Grassmann manifold are equivalence classes of $n \times p$, $p < n$ orthogonal matrices, where two matrices are equivalent if their columns span the same p -dimensional subspace.

In practice an element \mathbf{X} of $\mathcal{G}_{n,p}$ is represented by an orthonormal basis as a $n \times p$ matrix, i.e., $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$. The geodesic distance between two points on the Grassmann manifold can be computed as:

$$d_G(\mathbf{X}, \mathbf{Y}) = \|\Theta\|_2 \quad (5)$$

where $\Theta = [\theta_1, \theta_2, \dots, \theta_p]$ is the principal angle vector, i.e.:

$$\cos(\theta_k) = \max_{\mathbf{x} \in \mathbf{X}, \mathbf{y} \in \mathbf{Y}} \mathbf{x}^T \mathbf{y} = \mathbf{x}_k^T \mathbf{y}_k \quad (6)$$

$$\text{s.t.} : \quad \|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1$$

$$\mathbf{x}^T \mathbf{x}_i = 0; \quad i = 1, 2, \dots, k-1$$

$$\mathbf{y}^T \mathbf{y}_i = 0; \quad i = 1, 2, \dots, k-1$$

The principal angles have the property of $\theta_i \in [0, \pi/2]$ and can be computed through the SVD of $\widehat{\mathbf{X}}^T \mathbf{Y}$ (Edelman et al., 1999).

The prominent theme of analysing Grassmann manifolds is to embed them in higher dimensional Euclidean spaces. In this work, we are interested in embedding Grassmann manifolds in RKHS, which can be implicitly achieved through Grassmann kernels. A function $k : \mathcal{G}_{n,p} \times \mathcal{G}_{n,p} \rightarrow \mathbb{R}^+$ is a Grassmann kernel provided that it is positive definite and well defined for all $\mathbf{X} \in \mathcal{G}_{n,p}$. The well-defined criterion means that the kernel is invariant to various representations of the subspaces, i.e., $k(\mathbf{X}, \mathbf{Y}) = k(\mathbf{X}\mathbf{Q}_1, \mathbf{Y}\mathbf{Q}_2)$, $\forall \mathbf{Q}_1, \mathbf{Q}_2 \in \mathbb{O}(p)$, where $\mathbb{O}(p)$ indicates orthonormal matrices of order p (Hamm and Lee, 2008). The repertoire of Grassmann kernels includes Binet–Cauchy (Wolf and Shashua, 2003) and projection kernels (Hamm and Lee, 2008). Furthermore, the first canonical

³ Orthogonal group $\mathbb{O}(n)$ is the space of all $n \times n$ orthogonal matrices. It is not a vector space but a differentiable manifold with two connected components.

⁴ A relation \sim on manifold \mathcal{M} is an equivalence relation iff it is reflexive ($\mathbf{X} \sim \mathbf{X}, \forall \mathbf{X} \in \mathcal{M}$), symmetric ($\mathbf{X} \sim \mathbf{Y}$, iff $\mathbf{Y} \sim \mathbf{X}, \forall \mathbf{X}, \mathbf{Y} \in \mathcal{M}$) and transitive (if $\mathbf{X} \sim \mathbf{Y}$ and $\mathbf{Y} \sim \mathbf{Z}$ then $\mathbf{X} \sim \mathbf{Z}, \forall \mathbf{X}, \mathbf{Y}, \mathbf{Z} \in \mathcal{M}$). The set of all elements that are equivalent to a point \mathbf{X} is called the equivalence class of \mathbf{X} , i.e., $[\mathbf{X}] = \{\mathbf{Y} \in \mathcal{M} : \mathbf{Y} \sim \mathbf{X}\}$. The set $\mathcal{Y} = \mathcal{M}/\sim$ of all equivalence classes of \sim in \mathcal{M} i.e., $\mathcal{Y} = \{[\mathbf{X}] : \mathbf{X} \in \mathcal{M}\} = \{[\mathbf{Y} \in \mathcal{M} : \mathbf{Y} \sim \mathbf{X}] : \mathbf{X} \in \mathcal{M}\}$, is called the quotient of \mathcal{M} by \sim .

correlation of two subspaces forms a pseudo kernel⁵ on Grassmann manifolds (Harandi et al., 2011). The three kernels are shown below:

$$k_{\text{BC}}(\mathbf{X}, \mathbf{Y}) = \det(\mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X}) \quad (7)$$

$$k_{\text{proj}}(\mathbf{X}, \mathbf{Y}) = \text{tr}(\mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X}) \quad (8)$$

$$k_{\text{CC}}(\mathbf{X}, \mathbf{Y}) = \max_{\mathbf{x} \in \mathbf{X}, \mathbf{y} \in \mathbf{Y}} \mathbf{x}^T \mathbf{y} \quad (9)$$

4. Kernel analysis on Grassmann manifolds

Given a set of input/output data $\{(X_1, l_1), (X_2, l_2), \dots, (X_N, l_N)\}$, where $X_i \in \mathcal{G}_{n,p}$ is a Grassmann point and l_i is the corresponding class label from $\mathcal{L} = \{1, 2, \dots, C\}$, we are interested in optimisation problems in the form of Tikhonov regularisation (Tikhonov et al., 1977):

$$\max \{ \mathbf{J}(\langle \mathbb{W}, \Phi(X_1) \rangle, \dots, \langle \mathbb{W}, \Phi(X_N) \rangle, l_1, \dots, l_N) + \lambda \Omega(\mathbb{W}) : \mathbb{W} \in \mathcal{H} \} \quad (10)$$

Here, \mathcal{H} is a prescribed Hilbert space of dimension h (h could be infinity) equipped with an inner product $\langle \cdot, \cdot \rangle$, $\Omega : \mathcal{H} \rightarrow \mathbb{R}$ is a regulariser, and $\mathbf{J} : (\mathbb{R}^h)^N \times \mathcal{L}^N \rightarrow \mathbb{R}$ is a cost function. For certain choices of the regulariser, solving (10) reduces to identifying N parameters and not the dimension of \mathcal{H} . This is more formally explained by the representer theorem (Shawe-Taylor and Cristianini, 2004) which states that the solution $\widehat{\mathbb{W}}$ of (10) is a linear combination of the inputs when the regulariser is the square of the Hilbert space norm. For vector Hilbert spaces, this result is straightforward to prove and dates back to 1970s (Kimeldorf and Wahba, 1970). Argyriou et al. (2009) showed that the representer theorem holds for matrix Hilbert spaces as well.

In the following subsections, we first elucidate the details of the proposed Grassmann graph embedding DA (GGDA) algorithm, followed by how to use the mapping obtained by GGDA to accomplish classification. We then show that the conventional Grassmann DA (Hamm and Lee, 2008) is a special case of GGDA.

4.1. Grassmann graph embedding discriminant analysis (GGDA)

A graph (\mathbb{V}, \mathbb{G}) in our context refers to a collection of vertices or nodes, \mathbb{V} , and a collection of edges that connect pairs of vertices. We note that \mathbb{G} is a symmetric matrix with elements describing the similarity between pairs of vertices. Moreover, the diagonal matrix \mathbb{D} and the Laplacian matrix \mathbb{L} of a graph are defined as $\mathbb{L} = \mathbb{D} - \mathbb{G}$, with the diagonal elements of \mathbb{D} obtained as $\mathbb{D}(i, i) = \sum_j \mathbb{G}(i, j)$.

Given N labelled points $\mathbb{X} = \{(X_i, l_i)\}_{i=1}^N$ from the underlying Grassmann manifold $\mathcal{G}_{n,p}$, where $X_i \in \mathbb{R}^{n \times p}$ and $l_i \in \{1, 2, \dots, C\}$, with C denoting the number of classes, the local geometrical structure of $\mathcal{G}_{n,p}$ can be modelled by building a within-class similarity graph G_w and a between-class similarity graph G_b . The simplest forms of G_w and G_b are based on the nearest neighbour graphs defined below:

$$G_w(i, j) = \begin{cases} 1, & \text{if } X_i \in N_w(X_j) \text{ or } X_j \in N_w(X_i) \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

$$G_b(i, j) = \begin{cases} 1, & \text{if } X_i \in N_b(X_j) \text{ or } X_j \in N_b(X_i) \\ 0, & \text{otherwise} \end{cases} \quad (12)$$

In (11), $N_w(X_i)$ is the set of v_w neighbours $\{X_i^1, X_i^2, \dots, X_i^{v_w}\}$, sharing the same label as l_i . Similarly in (12), $N_b(X_i)$ contains v_b

⁵ A pseudo kernel is a function where the positive definiteness is not guaranteed to be satisfied for whole range of the function's parameters. Nevertheless, it is possible to convert a pseudo kernel into a true kernel, as discussed for example in (Chen et al., 2009).

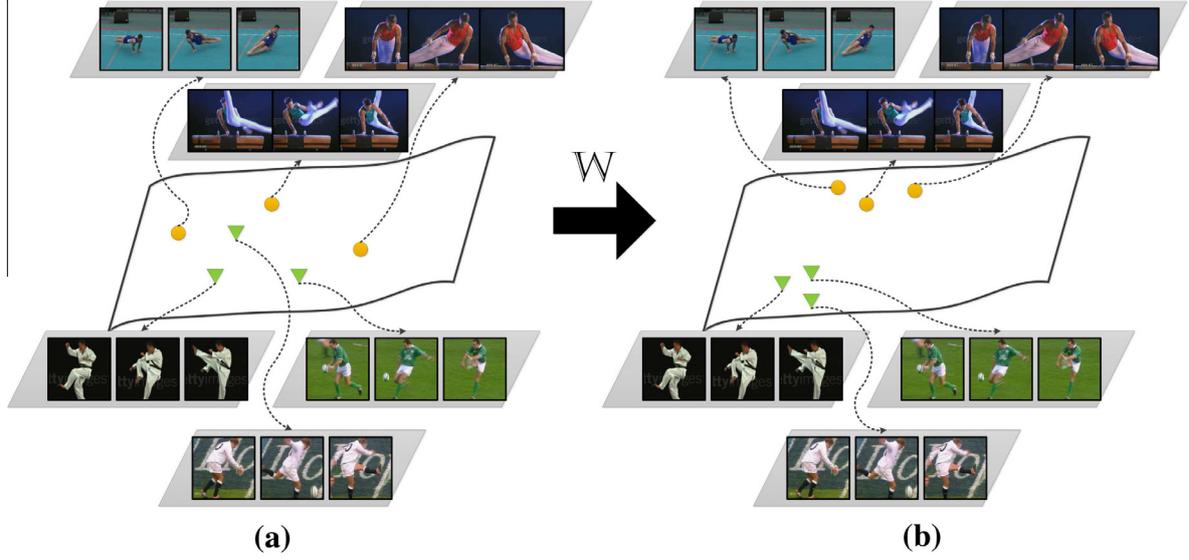


Fig. 2. A conceptual illustration of the proposed approach. (a) Actions can be modelled by linear subspaces, which in turn can be interpreted as points on a Grassmann manifold. In this figure, two types of actions (*kicking* and *swinging*) are shown. Having a proper geodesic distance between the points on the manifold, it is possible to convert the action recognition problem into a point to point classification problem. (b) Through the use of a Grassmann kernel, points on the Grassmann manifold can be mapped into an optimised RKHS, where not only certain local properties have been retained but also the discriminatory power between classes has been increased. Unlike conventional formalism of discriminant analysis, the proposed method preserves the geometrical structure and local information of a manifold by exploiting within-class and between-class similarity graphs.

neighbours having different labels. We note that more complex similarity graphs, like heat kernel graphs, can also be used to encode distances between points on Grassmann manifolds (Rosenberg, 1997).

Our aim is to simultaneously maximise a measure of discriminatory power and preserve the geometry of points (see Fig. 2 for a conceptual demonstration). This can be formalised by finding $\mathbb{W} : \Phi(X_i) \rightarrow Y_i$ such that the connected points of G_w are placed as close as possible, while the connected points of G_b are moved as far as possible. As such, a mapping must be sought by optimising the following two objective functions:

$$f_1 = \min \frac{1}{2} \sum_{ij} \|Y_i - Y_j\|^2 G_w(i, j) \quad (13)$$

$$f_2 = \max \frac{1}{2} \sum_{ij} \|Y_i - Y_j\|^2 G_b(i, j) \quad (14)$$

Eq. (13) punishes neighbours in the same class if they are mapped far away, while (14) punishes points of different classes if they are mapped close together. According to the representer theorem (Shawe-Taylor and Cristianini, 2004), the solution $\mathbb{W} = [\Gamma_1 | \Gamma_2 | \dots | \Gamma_r]$, can be expressed as a linear combination of data points, i.e., $\Gamma_i = \sum_{j=1}^N w_{ij} \phi(X_j)$. More specifically:

$$Y_i = (\langle \Gamma_1, \phi(X_i) \rangle, \langle \Gamma_2, \phi(X_i) \rangle, \dots, \langle \Gamma_r, \phi(X_i) \rangle)^T \quad (15)$$

We note that $\langle \Gamma_i, \phi(X_i) \rangle = \sum_{j=1}^N w_{ij} \text{tr}(\phi(X_j)^T \phi(X_i)) = \sum_{j=1}^N w_{ij} k(X_j, X_i)$, $Y_i = \mathbb{W}^T \mathbb{K}_i$, with $\mathbb{K}_i = (k(X_i, X_1), k(X_i, X_2), \dots, k(X_i, X_N))^T$ and

$$\mathbb{W} = \begin{pmatrix} w_{1,1} & w_{1,2} & \dots & w_{1,r} \\ w_{2,1} & w_{2,2} & \dots & w_{2,r} \\ \vdots & \vdots & \ddots & \vdots \\ w_{N,1} & w_{N,2} & \dots & w_{N,r} \end{pmatrix}$$

Plugging this into (13) results in:

$$\begin{aligned} \frac{1}{2} \sum_{ij} \|Y_i - Y_j\|^2 G_w(i, j) &= \frac{1}{2} \sum_{ij} \|\mathbb{W}^T \mathbb{K}_i - \mathbb{W}^T \mathbb{K}_j\|^2 G_w(i, j) \\ &= \sum_i \text{tr}(\mathbb{W}^T \mathbb{K}_i \mathbb{K}_i^T \mathbb{W}) G_w(i, i) \\ &\quad - \sum_{ij} \text{tr}(\mathbb{W}^T \mathbb{K}_j \mathbb{K}_i^T \mathbb{W}) G_w(i, j) \\ &= \text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{D}_w \mathbb{K}^T \mathbb{W}) - \text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{G}_w \mathbb{K}^T \mathbb{W}) \end{aligned} \quad (16)$$

where $\mathbb{K} = [\mathbb{K}_1 | \mathbb{K}_2 | \dots | \mathbb{K}_N]$. Considering that $\mathbb{L}_b = \mathbb{D}_b - \mathbb{W}_b$, in a similar manner it can be shown that (14) can be simplified to:

$$\begin{aligned} \frac{1}{2} \sum_{ij} \|Y_i - Y_j\|^2 G_b(i, j) &= \text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{D}_b \mathbb{K}^T \mathbb{W}) - \text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{G}_b \mathbb{K}^T \mathbb{W}) \\ &= \text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{L}_b \mathbb{K}^T \mathbb{W}) \end{aligned} \quad (17)$$

To solve (13) and (14) simultaneously, we need to add the following normalising constraint to the problem:

$$\text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{D}_w \mathbb{K}^T \mathbb{W}) = 1 \quad (18)$$

This constraint enables us to convert the minimisation problem (13) into a maximisation one. Consequently, both equations can be combined into one maximisation problem. Moreover, as shown later, the imposed constraint acts as a norm regulariser in the original Tikhonov problem (10), thus satisfying the necessary condition of the representer theorem. Plugging (18) into (13) results in:

$$\begin{aligned} \min \{ \text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{D}_w \mathbb{K}^T \mathbb{W}) - \text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{G}_w \mathbb{K}^T \mathbb{W}) \} \\ = \min \{ 1 - \text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{G}_w \mathbb{K}^T \mathbb{W}) \} = \max \{ \text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{G}_w \mathbb{K}^T \mathbb{W}) \} \end{aligned} \quad (19)$$

subject to the constraint shown in (18). As a result, the max versions of (13) and (14) can be merged by the Lagrangian method (C.M. Bishop, 2006) as follows:

$$\begin{aligned} \max \{ \text{tr}(\mathbb{W}^T \mathbb{K} (\mathbb{L}_b + \beta \mathbb{G}_w) \mathbb{K}^T \mathbb{W}) \} \\ \text{subject to } \text{tr}(\mathbb{W}^T \mathbb{K} \mathbb{D}_w \mathbb{K}^T \mathbb{W}) = 1 \end{aligned} \quad (20)$$

where β is a Lagrangian multiplier. The solution to the optimisation in (20) can be sought as the r largest eigenvectors of the following generalised eigenvalue problem:

$$\mathbb{K}\{\mathbb{L}_b + \beta\mathbb{G}_w\}\mathbb{K}^T\mathbb{W} = \lambda\mathbb{K}\mathbb{D}_w\mathbb{K}^T\mathbb{W} \quad (21)$$

We note that in (21), the imposed constraint (18) serves as a norm regulariser and satisfies the representer theorem condition. Algorithm 1 assembles all the above details into pseudo-code for training the Grassmann graph embedding discriminant analysis (GGDA).

Algorithm 1. Pseudocode for training Grassmann graph-embedding discriminant analysis (GGDA)

Input: Training set $\mathbb{X} = \{(X_i, l_i)\}_{i=1}^N$ from the underlying Grassmann manifold $\mathcal{G}n, p$, where $X_i \in \mathbb{R}^{n \times p}$ is a subspace and $l_i \in \{1, 2, \dots, C\}$, with C denoting the number of classes

- A kernel function k_{ij} , for measuring the similarity between two points on a Grassmann manifold

Output: The projection matrix $\mathbb{W} = [\mathbf{Y}_1 | \mathbf{Y}_2 | \dots | \mathbf{Y}_r]$,

- 1: Compute the Gram matrix $[\mathbb{K}]_{ij}$ for all X_i, X_j
- 2: **for** $i = 1 \rightarrow N - 1$ **do**
- 3: **for** $j = i + 1 \rightarrow N$ **do**
- 4: Compute the geodesic distances $d_g(i, j)$ between X_i and X_j .
- 5: $d_g(j, i) = d_g(i, j)$
- 6: **end for**
- 7: **end for**
- 8: $\mathbb{G}_w \leftarrow \mathbf{0}_{N \times N}$
- 9: $\mathbb{G}_b \leftarrow \mathbf{0}_{N \times N}$
- % Use the obtained $d_g(i, j)$ to determine neighbourhoods in the following loop.
- 10: **for** $i = 1 \rightarrow N$ **do**
- 11: **if** (X_j is in the first k_w nearest neighbours of X_i) **and** ($l_j == l_i$) **then**
- 12: $G_w(i, j) \leftarrow 1$
- 13: $G_w(j, i) \leftarrow 1$
- 14: **end if**
- 15: **if** (X_j is in the first k_b nearest neighbours of X_i) **and** ($l_j \neq l_i$) **then**
- 16: $G_b(i, j) \leftarrow 1$
- 17: $G_b(j, i) \leftarrow 1$
- 18: **end if**
- 19: **end for**
- 20: $\mathbb{D}_w \leftarrow \mathbf{0}_{N \times N}$
- 21: $\mathbb{D}_b \leftarrow \mathbf{0}_{N \times N}$
- 22: $D_w(i, i) \leftarrow \sum_j G_w(i, j)$
- 23: $D_b(i, i) \leftarrow \sum_j G_b(i, j)$
- 24: $\mathbb{L}_b \leftarrow \mathbb{D}_b - \mathbb{G}_b$
- 25: $\{\mathbf{Y}_i, \tilde{\lambda}_i\}_{i=1}^r \leftarrow$ generalised eigenvectors and eigenvalues of $\mathbb{K}\{\mathbb{L}_b + \beta\mathbb{G}_w\}\mathbb{K}^T\mathbb{W} = \lambda\mathbb{K}\mathbb{D}_w\mathbb{K}^T\mathbb{W}$, with the eigenvectors ordered according to descending eigenvalues.

4.2. Classification

Upon acquiring the mapping \mathbb{W} , classification tasks on Grassmann manifolds are reduced to classification tasks in vector spaces. More precisely, for any query image set X_q , a vector representation using the kernel function and the mapping \mathbb{W} is acquired, i.e., $\mathbf{v}_q = \mathbb{W}^T \mathbb{K}_q$, where $\mathbb{K}_q = (\langle \phi(X_1), \phi(X_q) \rangle, \langle \phi(X_2), \phi(X_q) \rangle, \dots, \langle \phi(X_N), \phi(X_q) \rangle)^T$. Similarly, gallery points X_i are represented by r dimensional vectors $\mathbf{v}_i = \mathbb{W}^T \mathbb{K}_i$. Classification methods such as nearest-neighbours or support vector machines (C.M. Bishop, 2006) can be employed to label \mathbf{v}_q .

4.3. Relation to Grassmann discriminant analysis

Here we address the relation between GGDA and Grassmann discriminant analysis (GDA) (Hamm and Lee, 2008). GDA is a learn-

ing framework on Grassmann manifolds that utilises kernel analysis to project Grassmann points into a higher discriminative Hilbert space. More specifically, let $\mathbb{X} = \{(X_i, l_i)\}_{i=1}^N$ be a set of N labelled points on Grassmann manifold $\mathcal{G}n, p$. GDA seeks a transformation $\mathbb{W} : \Phi(X_i) \rightarrow Y_i$ such that the ratio of between-class to within-class scatters is maximised. The within-class and between-class scatters are defined as:

$$S_W = \sum_{j=1}^C \sum_{i: l_i=j} \|Y_i - \mu_j\|^2 \quad (22)$$

$$S_B = \sum_{j=1}^C n_j \|\bar{\mu} - \mu_j\|^2 \quad (23)$$

where $\sum_{i: l_i=j}$ denotes the summation over i such that $l_i = j$, while $\mu_j = \frac{1}{n_j} \sum_{i: l_i=j} Y_i$ is the mean of the samples in class j , and $\bar{\mu} = \frac{1}{N} \sum_{j=1}^C n_j \mu_j$ is the mean of all samples.

We note that if the data pairs in the same class are moved closer, the within-class scatter S_W gets smaller. On a similar note, if the data pairs in different classes are more separated from each other, the between-class scatter S_B gets larger. GDA can be seen as a special case of GGDA when particular within-class and between-class similarity graphs are used. The following lemma formalises the relation between GGDA and GDA.

Lemma 1. GDA is a special case of GGDA if

$$G_w(i, j) = \begin{cases} \frac{1}{n_k}, & \text{if } l_i = l_j = k \\ 0, & \text{if } l_i \neq l_j \end{cases} \quad (24)$$

$$G_b(i, j) = \begin{cases} \frac{1}{N} - \frac{1}{n_k}, & \text{if } l_i = l_j = k \\ \frac{1}{N}, & \text{if } l_i \neq l_j \end{cases} \quad (25)$$

A proof of this lemma is given in A.

5. Experiments

In this section we first compare and contrast the performance of the proposed GGDA method against several state-of-the-art approaches on the UCF sport action dataset (Rodriguez et al., 2008), the KTH human motion dataset (Schuldts et al., 2004) and the Ballet dataset (Wang and Mori, 2009). We then conclude the section by assessing the sensitivity of the proposed method against occlusion and misalignment.

5.1. Empirical evaluations

Here, we appraise the performance of the proposed GGDA method⁶ against the state-of-the-art Grassmann discriminant analysis (Hamm and Lee, 2008), kernel version of affine hull image-set distance (Cevikalp and Triggs, 2010), tensor canonical correlation analysis (Kim and Cipolla, 2009), spatial-temporal words (Niebles et al., 2008) and hierarchy of discriminative space-time neighbourhood features (Kovashka and Grauman, 2010) on the UCF sport action dataset (Rodriguez et al., 2008), the KTH human motion dataset (Schuldts et al., 2004) and the Ballet dataset (Wang and Mori, 2009). In all experiments, the projection kernel has been used in GGDA.

5.1.1. UCF SPORT dataset

The UCF sport action dataset (Rodriguez et al., 2008) consists of ten categories of human actions including swinging on the pommel

⁶ Source code for the proposed method is available at <<http://itee.uq.edu.au/uqmhara1>>.



Fig. 3. Examples from: (a) UCF sport action dataset (Rodriguez et al., 2008), (b) KTH dataset (Schuldt et al., 2004), and (c) Ballet dataset (Wang and Mori, 2009).

Table 4

Recognition accuracy (in %) for the KTH action recognition dataset using spatio-temporal words (STW) (Niebles et al., 2008), fusion of appearance and distribution method (BoW with MKL) (Bregonzio et al., 2012), tensor canonical correlation analysis (TCCA) (Kim and Cipolla, 2009), Grassmann discriminant analysis (GDA) (Hamm and Lee, 2008) and the proposed GGDA approach.

Method	Recognition accuracy (%)
STW (Niebles et al., 2008)	83
BoW-MKL (Bregonzio et al., 2012)	94
TCCA (Kim and Cipolla, 2009)	95
GDA (Hamm and Lee, 2008) with image-set modelling	83
GDA (Hamm and Lee, 2008) with ARMA modelling	86
GGDA with image-set modelling	97
GGDA with ARMA modelling	99

(LOO) cross validation protocol used in (Niebles et al., 2008; Kim and Cipolla, 2009). The classification results are reported in Table 4.

We compared GGDA against two state-of-the-art Euclidean approaches: spatial-temporal words (STW) (Niebles et al., 2008) and bag of words model in conjunction with multiple kernel learning (Bach et al., 2004) (BoW-MKL) (Bregonzio et al., 2012). In STW, a video sequence is represented by a set of spatial-temporal words, extracted from space-time interest points. The algorithm then utilises latent topic models such as the probabilistic latent semantic

Table 6

Recognition accuracy (in %) along its standard deviation for the Ballet dataset using Grassmann geodesic distance (Turaga et al., 2011), Kernel Affine Hull method (KAHM) (Cevikalp and Triggs, 2010), Grassmann discriminant analysis (GDA) (Hamm and Lee, 2008) and the proposed GGDA approach.

Method	Recognition accuracy
Geodesic distance	77.34% \pm 1.8
KAHM (Cevikalp and Triggs, 2010)	79.71% \pm 2.3
GDA (Hamm and Lee, 2008) with image-set modelling	78.05% \pm 2.9
GDA (Hamm and Lee, 2008) with ARMA modelling	73.70% \pm 1.8
GGDA with image-set modelling	83.08% \pm 1.8
GGDA with ARMA modelling	77.63% \pm 2.9

analysis (Hofmann, 1999) to learn the probability distributions of the spatial-temporal words. BoW-MKL exploits global spatio-temporal distribution of interest points by extracting holistic features from clouds of interest points accumulated over multiple temporal scales. Then extracted features are fused using MKL. We also compared GGDA against Tensor Canonical Correlation Analysis (TCCA) (Kim and Cipolla, 2009) and conventional discriminant analysis on Grassmann manifolds (GDA) (Hamm and Lee, 2008). TCCA is an extension of canonical correlation analysis (a principled tool to inspect linear relations between two sets of vectors) to tensor spaces and measures video-to-video volume similarity.

Table 5

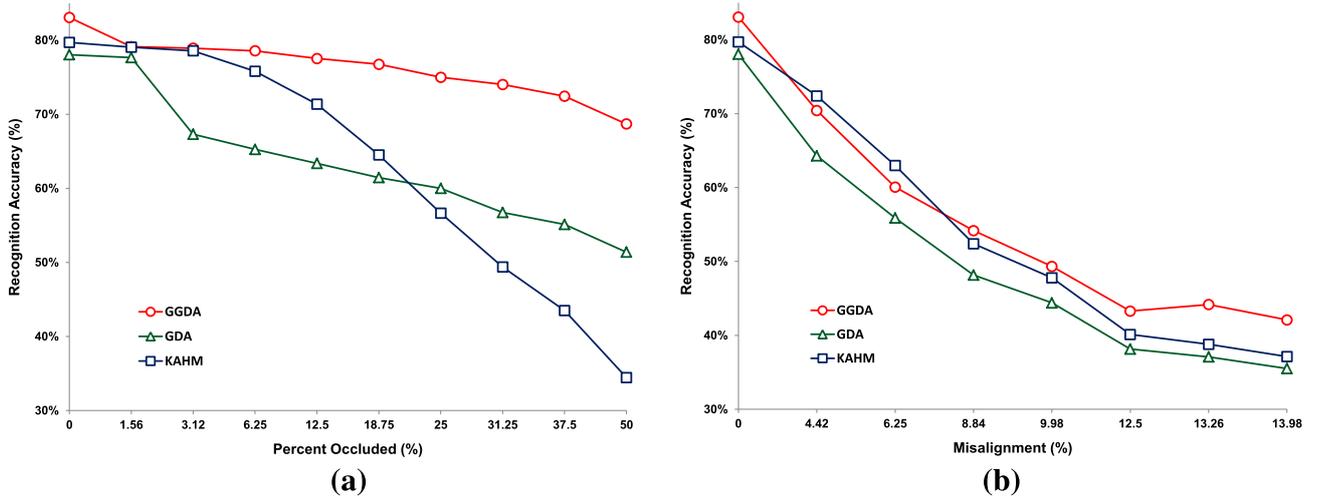
Confusion matrix (in %) for the GGDA method on the KTH action recognition dataset using LOO protocol.

	Boxing	Hand clapping	Hand waving	Jogging	Running	Walking
Boxing	100	0	0	0	0	0
Hand clapping	0	100	0	0	0	0
Hand waving	0	0	100	0	0	0
Jogging	0	0	0	99	0	1
Running	0	0	0	1	97	2
Walking	0	0	0	0	2	98

Table 7

Confusion matrix (in %) for the GGDA method on the Ballet dataset. Actions are modelled by image-sets.

	LR hand opening	RL hand opening	Standing hand opening	Leg swinging	Jumping	Turning	Hopping	Standing still
LR hand opening	81.50	2.50	2.50	0.00	3.00	5.00	0.00	5.50
RL hand opening	0.00	100.00	0.00	0.00	0.00	0.00	0.00	0.00
Standing hand opening	0.00	0.00	87.88	4.55	0.00	6.06	0.00	1.52
Leg swinging	3.57	0.00	0.00	85.71	3.57	3.57	0.00	3.57
Jumping	4.55	4.55	9.09	13.14	50.54	0.00	13.60	4.55
Turning	0.00	0.00	5.82	0.00	0.00	94.18	0.00	0.00
Hopping	0.00	2.33	17.23	2.33	2.33	0.00	71.14	4.65
Standing still	6.00	0.00	27.50	11.00	0.00	10.00	0.00	45.50

**Fig. 4.** Recognition rate on the Ballet dataset (Wang and Mori, 2009) for increasing amount of (a) occlusion and (b) misalignment, for KAHM (Cevikalp and Triggs, 2010), GDA (Hamm and Lee, 2008) and the proposed GGDA algorithm.

Looking at the results in Table 4, the proposed GGDA method achieves the highest classification accuracy. We note that though walking, jogging and running are sometimes confused by GGDA (see the confusion matrix in Table 5). However, in comparison to the other methods GGDA is able to reduce this ambiguity greatly.

5.1.3. Ballet dataset

The previous experiments may imply that actions should be exclusively modelled by an ARMA process. As we show in this experiment, this is not always the case, indicating that image-sets are also useful for modelling of actions. The Ballet dataset contains 44 real video sequences of 8 actions collected from an instructional ballet DVD (Wang and Mori, 2009).⁷ The dataset consists of 8 complex motion patterns performed by three subjects. The actions include: ‘left-to-right hand opening’, ‘right-to-left hand opening’, ‘standing hand opening’, ‘leg swinging’, ‘jumping’, ‘turning’, ‘hopping’ and ‘standing still’. Fig. 3(c) shows samples. This dataset is very challenging due to the significant intra-class variations in terms of speed, spatial and temporal scale, clothing and movement.

Available samples of each action were randomly split into training and testing set (the number of actions in both training and testing sets were fairly even). The process of random splitting was repeated ten times and the average classification accuracy is reported in Table 6. For comparison, the GGDA algorithm is contrasted with geodesic distance on Grassmann manifolds (Eq. (5)), and the state-of-art kernel version of affine hull set matching (KAHM) (Cevikalp and Triggs, 2010) and Grassmann discriminant analysis (GDA) (Hamm and Lee, 2008). Cevikalp and Triggs

(2010) proposed to measure the similarity between image-sets using geometric distances between their convex models. In this experiment, for Grassmann-based analysis, actions were modelled by image-sets of order 10. For the sake of comparison, we also modelled actions by ARMA process with state-space dimension $p = 20$ (the observability matrix was truncated at 5).

Table 6 shows that the GGDA algorithm with image-set modelling obtains the highest accuracy and outperforms state-of-the-art methods KAHM and GDA significantly. The confusion matrix of the proposed GGDA method is shown in Table 7. We note that the lowest recognition accuracy belongs to the ‘standing still’ action which is mainly confused with ‘standing hand opening’. We conjecture that the low performance of GGDA method with ARMA modelling is due to lack of temporal information for some action, such as the ‘standing still’ action.

5.2. Sensitivity analysis

The performance of a visual recognition system can be negatively affected by variations in the environment (e.g., illumination), capturing device (e.g., pose variation, occlusion) as well as preprocessing steps that prepare data for the system (e.g., registration). A detailed sensitivity analysis for the aforementioned factors is beyond the scope of this paper. However, in this paper we analyse the sensitivity of GGDA algorithm against occlusion and misalignment. To this end, we elect the Ballet dataset (Wang and Mori, 2009) for our analysis since this dataset has a uniform background and fair illumination (and therefore minimises the effect of variations in illumination and background in the analysis). Actions are modelled by image-sets and the test setup employed in Section 5.1.3 is utilised again.

⁷ The study in (Wang and Mori, 2009) addresses the problem of recognising actions in still images, which is different from the work presented here.

5.2.1. Sensitivity to occlusion

In this experiment, we assess the performance at various levels of occlusion, from 1.56% up to 50%, by replacing a set of randomly located square blocks of size 4×4 in a test image with a blank block. The location of occlusion is randomly chosen for each test image and is unknown to the system. The training images do not contain occlusions. Methods that select fixed features or blocks of the image are unlikely to succeed here due to the unpredictable location of the occlusion.

Fig. 4(a) shows the recognition rates of KAHM, GDA and GGDA. The proposed GGDA method significantly outperforms the other two methods for almost all levels of occlusion. Up to 40 percent occlusion, the performance of GGDA has dropped roughly by 10 percentage points. While robustness to occlusion can be partially attributed to subspace modelling (as can also be seen for GDA), we conjecture that the proposed GGDA method has better captured the true Grassmannian geometry (through within and between graphs as the training images do not contain occlusions) and is hence more robust to the missing parts.

5.2.2. Sensitivity to misalignment

The temporal and spatial misalignment could deteriorate the performance of an action recogniser drastically (Shariat and Pavlovic, 2011). In this work we only consider spatial misalignment and assess and contrast the sensitivity of GGDA algorithm as compared to KAHM (Cevikalp and Triggs, 2010) and GDA (Hamm and Lee, 2008) methods. To this end, we have introduced random displacements to the frames of query videos and measured the accuracy for various amounts of displacements. Fig. 4(b) shows the results of misalignment analysis for KAHM (Cevikalp and Triggs, 2010), GDA (Hamm and Lee, 2008) and the proposed GGDA methods. The horizontal axis here demonstrates the degree of misalignment, i.e., the length of random displacement vector divided by the maximum possible misalignment (for frame of size $S_x \times S_y$, the maximum possible misalignment is $\frac{1}{2}\sqrt{S_x^2 + S_y^2}$). Fig. 4(b) reveals that all studied algorithms are sensitive to misalignment. The larger the displacement, the lower would be the recognition accuracy. This is mainly due to the holistic representation of images which has been shown to be highly fragile to misalignment (Wong et al., 2012).

6. Main findings and future directions

In this paper, we first demonstrated how actions can be modelled by linear subspaces. Subspaces are able to accommodate the effects of various image variations and can capture the dynamic properties of videos (Turaga et al., 2011). Since subspaces form non-Euclidean and curved Riemannian manifolds known as Grassmann manifolds, we exploited the geometry of space to design an action recogniser. As such, we proposed graph-embedding discriminant analysis on Grassmann manifolds by embedding manifolds into RKHS. The proposed method utilises within-class and between-class similarity graphs to characterise intra-class compactness and inter-class separability, respectively.

Thorough experiments on the KTH (Schuldt et al., 2004), UCF Sports (Rodriguez et al., 2008) and Ballet (Wang and Mori, 2009) datasets, which include various realistic challenges such as background clutter, partial occlusion, changes in viewpoint, scale and illumination, and complexity of motion showed that the proposed approach obtains notable improvements in discrimination accuracy in comparison to several state-of-the-art methods. This included Grassmann discriminant analysis (Hamm and Lee, 2008), kernel version of affine hull image-set distance (Cevikalp and Triggs, 2010), tensor canonical correlation analysis (Kim and Cipolla, 2009), spatial-temporal words (Niebles et al., 2008) and

hierarchy of discriminative space-time neighbourhood features (Kovashka and Grauman, 2010).

The proposed GGDA algorithm (like other discriminant analysis techniques) requires several samples of each class to determine the optimum mapping. As such, GGDA cannot be applied to action recognition problems where just one sample video per class is available for training. Moreover, since the structure of manifold is encoded via between and within similarity graphs, inappropriate parameter selection (for example the size of neighbourhood N_w in Eqn. 11) might result in poor performance.

In this paper, each action is modelled by just one subspace. When the length of action video is small, this straightforward treatment appears to be appropriate. However, for very complex motions or when an action is described by a lengthy video, straightforward subspace modelling might not be adequate enough. As such, future avenues of research include exploring how several subspaces can be generated from complex motion videos. We are also keen to explore how discriminatory subspaces can be generated from interest-points (e.g., local spatio-temporal features) in a video.

Acknowledgements

NICTA is funded by the Australian Government as represented by the *Department of Broadband, Communications and the Digital Economy*, as well as the Australian Research Council through the *ICT Centre of Excellence* program.

Appendix A. Proof of Lemma 1

Plugging (24) into (13), we get:

$$\begin{aligned}
\frac{1}{2} \sum_{ij} \|\mathbf{Y}_i - \mathbf{Y}_j\|^2 G_w(i, j) &= \sum_{ij} \mathbf{Y}_i^T \mathbf{Y}_j G_w(i, j) - \sum_{ij} \mathbf{Y}_i^T \mathbf{Y}_j G_w(i, j) \\
&= \sum_i \mathbf{Y}_i^T \mathbf{Y}_i \sum_j G_w(i, j) - \sum_{k=1}^c \frac{1}{n_k} \sum_{i,j:l_i=l_j=k} \mathbf{Y}_i^T \mathbf{Y}_j \\
&= \sum_i \mathbf{Y}_i^T \mathbf{Y}_i - \sum_{k=1}^c \frac{1}{n_k} \sum_{i,j:l_i=l_j=k} \mathbf{Y}_i^T \mathbf{Y}_j \\
&= \sum_{k=1}^c \sum_{i:l_i=k} \left(\mathbf{Y}_i - \frac{1}{n_k} \sum_{j:l_j=k} \mathbf{Y}_j \right)^T \left(\mathbf{Y}_i - \frac{1}{n_k} \sum_{j:l_j=k} \mathbf{Y}_j \right) \\
&= \sum_{k=1}^c \sum_{i:l_i=k} (\mathbf{Y}_i - \boldsymbol{\mu}_k)^2 = S_w \tag{A.1}
\end{aligned}$$

To show the equivalency for between-class scatters, we note that:

$$S_M = S_B + S_W = \sum_{i=1}^N \|\mathbf{Y}_i - \boldsymbol{\mu}\|^2 = \sum_{i=1}^N \mathbf{Y}_i^T \mathbf{Y}_i - \frac{1}{N} \sum_{i,j=1}^N \mathbf{Y}_i^T \mathbf{Y}_j \tag{A.2}$$

By plugging (25) into (14) and using $G_b(i, j) = \frac{1}{N} - G_w(i, j)$, we get:

$$\begin{aligned}
\frac{1}{2} \sum_{ij} \|\mathbf{Y}_i - \mathbf{Y}_j\|^2 G_b(i, j) &= \frac{1}{2N} \sum_{ij} \|\mathbf{Y}_i - \mathbf{Y}_j\|^2 - \frac{1}{2} \sum_{ij} \|\mathbf{Y}_i \\
&\quad - \mathbf{Y}_j\|^2 G_w(i, j) \\
&= \frac{1}{N} \sum_{ij} \mathbf{Y}_i^T \mathbf{Y}_i - \frac{1}{N} \sum_{ij} \mathbf{Y}_i^T \mathbf{Y}_j - S_w \\
&= \sum_i \mathbf{Y}_i^T \mathbf{Y}_i - \frac{1}{N} \sum_{ij} \mathbf{Y}_i^T \mathbf{Y}_j - S_w \\
&= S_M - S_W = S_B \tag{A.3}
\end{aligned}$$

References

- Argyriou, A., Micchelli, C.A., Pontil, M., 2009. When is there a representer theorem? Vector versus matrix regularizers. *J. Machine Learn. Res.* 10, 2507–2529.
- Bach, F.R., Lanckriet, G.R.G., Jordan, M.I., 2004. Multiple kernel learning, conic duality, and the SMO algorithm. In: *Internat. Conf. on Machine Learning*. ACM, New York, NY, USA.
- Bishop, C.M., 2006. *Pattern Recognition and Machine Learning*. Springer.
- Bregonzio, M., Xiang, T., Gong, S., 2012. Fusing appearance and distribution information of interest points for action recognition. *Pattern Recognition* 45 (3), 1220–1234.
- Cevikalp, H., Triggs, B., 2010. Face recognition based on image sets. In: *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp. 2567–2573.
- Chen, J., Ye, J., Li, Q., 2007. Integrating global and local structures: A least squares framework for dimensionality reduction. In: *Proc. Computer Vision and Pattern Recognition (CVPR)*, pp. 1–8.
- Chen, Y., Garcia, E.K., Gupta, M.R., Rahimi, A., Cazzanti, L., 2009. Similarity-based classification: Concepts and algorithms. *J. Machine Learn. Res.* 10, 747–776.
- Cock, K.D., Moor, B.D., 2002. Subspace angles between ARMA models. *Syst. Control Lett.* 46, 265–270.
- Comaniciu, D., Meer, P., 2002. Mean shift: A robust approach toward feature space analysis. *IEEE Trans. Pattern Anal. Machine Intell.* 24 (5), 603–619.
- Dalal, N., Triggs, W., 2005. Histograms of oriented gradients for human detection. In: *IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 886–893.
- Edelman, A., Arias, T.A., Smith, S.T., 1999. The geometry of algorithms with orthogonality constraints. *SIAM J. Matrix Anal. Appl.* 20 (2), 303–353.
- Hamm, J., Lee, D.D., 2008. Grassmann discriminant analysis: A unifying view on subspace-based learning. In: *Internat. Conf. on Machine Learning (ICML)*, pp. 376–383.
- Harandi, M.T., Sanderson, C., Shirazi, S., Lovell, B.C., 2011. Graph embedding discriminant analysis on Grassmannian manifolds for improved image set matching. In: *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp. 2705–2712.
- Harandi, M.T., Sanderson, C., Wiliem, A., Lovell, B.C., 2012. Kernel analysis over Riemannian manifolds for visual recognition of actions, pedestrians and textures. In: *IEEE Workshop on Applications of Computer Vision (WACV)*, pp. 433–439.
- Hofmann, T., 1999. Probabilistic latent semantic indexing. In: *Internat. ACM SIGIR Conf. on Research and Development in Information Retrieval*, pp. 50–57.
- Kim, T.-K., Cipolla, R., 2009. Canonical correlation analysis of video volume tensors for action categorization and detection. *IEEE Trans. Pattern Anal. Machine Intell. (PAMI)* 31, 1415–1428.
- Kimeldorf, G.S., Wahba, G., 1970. A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. *Ann. Math. Stat.* 41, 495–502.
- Kovashka, A., Grauman, K., 2010. Learning a hierarchy of discriminative space-time neighborhood features for human action recognition. In: *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp. 2046–2053.
- Li, B., Ayazoglu, M., Mao, T., Camps, O., Sznai, M., 2011. Activity recognition using dynamic subspace angles. In: *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp. 3193–3200.
- Lui, Y.M., 2012. Advances in matrix manifolds for computer vision. *Image Vision Comput.* 30 (6–7), 380–388.
- Niebles, J., Wang, H., Fei-Fei, L., 2008. Unsupervised learning of human action categories using spatial-temporal words. *Internat. J. Comput. Vision* 79, 299–318.
- O'Hara, S., Lui, Y.M., Draper, B.A., 2012. Using a product manifold distance for unsupervised action recognition. *Image Vision Comput.* 30 (3), 206–216.
- Rodriguez, M., Ahmed, J., Shah, M., 2008. Action MACH a spatio-temporal maximum average correlation height filter for action recognition. In: *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp. 1–8.
- Rosenberg, S., 1997. *The Laplacian on a Riemannian Manifold: An Introduction to Analysis on Manifolds*. Cambridge University Press.
- Schuld, C., Laptev, I., Caputo, B., 2004. Recognizing human actions: A local SVM approach. In: *Internat. Conf. on Pattern Recognition (ICPR)*, pp. 32–36.
- Shariat, S., Pavlovic, V., 2011. Isotonic CCA for sequence alignment and activity recognition. In: *IEEE Internat. Conf. on Computer Vision (ICCV)*. IEEE, pp. 2572–2578.
- Shawe-Taylor, J., Cristianini, N., 2004. *Kernel Methods for Pattern Analysis*. Cambridge University Press.
- Shirazi, S., Harandi, M.T., Sanderson, C., Lovell, B.C., 2012. Clustering on Grassmann manifolds via kernel embedding with application to action analysis. In: *IEEE Internat. Conf. on Image Processing*.
- Subbarao, R., Meer, P., 2009. Nonlinear mean shift over Riemannian manifolds. *Int. J. Comput. Vision* 84 (1), 1–20.
- Tikhonov, A.N., Arsenin, V.Y., 1977. *Solutions of Ill-Posed Problems*. V.H. Winston & Sons/John Wiley & Sons, Washington, D.C./New York.
- Turaga, P., Chellappa, R., 2009. Locally time-invariant models of human activities using trajectories on the Grassmannian. In: *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp. 2435–2441.
- Turaga, P., Chellappa, R., Subrahmanian, V., Udrea, O., 2008. Machine recognition of human activities: A survey. *IEEE Trans. Circ. Systems Video Technol.* 18 (11), 1473–1488.
- Turaga, P., Veeraraghavan, A., Srivastava, A., Chellappa, R., 2011. Statistical computations on Grassmann and Stiefel manifolds for image and video-based recognition. *IEEE Trans. Pattern Anal. Machine Intell.* 33 (11), 2273–2286.
- Tuzel, O., Porikli, F., Meer, P., 2008. Pedestrian detection via classification on Riemannian manifolds. *IEEE Trans. Pattern Anal. Machine Intell. (PAMI)* 30, 1713–1727.
- Wang, H., Ullah, M.M., Kläser, A., Laptev, I., Schmid, C., 2009. Evaluation of local spatio-temporal features for action recognition. In: *British Machine Vision Conference (BMVC)*.
- Wang, T., Shi, P., 2009. Kernel Grassmannian distances and discriminant analysis for face recognition from image sets. *Pattern Recognition Lett.* 30 (13), 1161–1165.
- Wang, Y., Mori, G., 2009. Human action recognition by semilattice topic models. *IEEE Trans. Pattern Anal. Machine Intell.* 31 (10), 1762–1774.
- Weinland, D., Ronfard, R., Boyer, E., 2011. A survey of vision-based methods for action representation, segmentation and recognition. *Computer Vision and Image Understanding* 115 (2), 224–241.
- Wolf, L., Shashua, A., 2003. Learning over sets using kernel principal angles. *J. Machine Learn. Res.* 4, 913–931.
- Wong, Y., Harandi, M., Sanderson, C., Lovell, B.C., 2012. On robust biometric identity verification via sparse encoding of faces: Holistic vs local approaches. In: *Internat. Joint Conf. on Neural Networks (IJCNN)*, pp. 1–8.
- Wu, X., Xu, D., Duan, L., Luo, J., 2011. Action recognition using context and appearance distribution features. In: *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp. 489–496.
- Yan, S., Xu, D., Zhang, B., Zhang, H.-J., Yang, Q., Lin, S., 2007. Graph embedding and extensions: A general framework for dimensionality reduction. *IEEE Trans. Pattern Anal. Machine Intell.* 29 (1), 40–51.