

Discriminative Common Vectors with Kernels

Review

Miguel A. Vezanzones

Grupo de Inteligencia Computacional
Universidad del País Vasco

Outline

- 1 Introduction
 - Background
 - FLD's modifications
- 2 Discriminative Common Vector (DCV)
 - Introduction
 - DCV by using the Range Space of S_W
 - DCV by using Difference Subspaces and the Gram-Schmidt Orthogonalization Procedure
- 3 Rough DCV
- 4 Kernel DCV

Outline

- 1 Introduction
 - Background
 - FLD's modifications
- 2 Discriminative Common Vector (DCV)
 - Introduction
 - DCV by using the Range Space of S_W
 - DCV by using Difference Subspaces and the Gram-Schmidt Orthogonalization Procedure
- 3 Rough DCV
- 4 Kernel DCV

Context

- Face recognition problem.
- Each image is represented by a vector in a wh -dimensional space.
- This space is called the sample space or the image space, and its dimension is typically very high.
- There is redundant information.

Objective

- Find a subspace based features extraction method that could succeed dealing with the small sample size problem.
- Small sample size problem: when data space dimensionality is larger than the number of samples in the training set.
- Subspace based methods:
 - Principal Component Analysis (PCA) -> unsupervised
 - Independent Component Analysis (ICA) -> unsupervised
 - Fisher's Linear Discriminant (FLD) -> supervised

Principal Component Analysis (PCA)

- Projection that maximizes the norm of the transformed total scatter matrix (covariances).

$$J_{PCA}(W_{opt}) = \arg \max_W |W^T S_T W|$$

- The optimal transformation is given by the eigenvectors of S_T .
- The PCA's projection directions are also called eigenfaces. Any face image in the sample space can be approximated by a linear combination of the significant eigenfaces.
- Tends to model unwanted within-class variations (lighting, expressions, occlusions,...) and the resulting classes tend to have more overlapping than other approaches.

Fisher's Linear Discriminant (FLD)

- Overcomes the limitations of the Eigenfaces method.
- Projections that maximize the between class scatter matrix, S_B , and minimize the within class scatter matrix, S_W .

$$J_{FLD}(W_{opt}) = \arg \max_W \frac{|W^T S_B W|}{|W^T S_W W|}$$

- The maximum is given by the eigenvectors of $S_W^{-1} S_B$.
- Not applicable within “small sample size problem” because S_W is singular in this case.

Scatter matrices

- Between class scatter matrix:

$$S_B = \sum_{i=1}^C N_i (\mu_i - \mu)(\mu_i - \mu)^T$$

- Within class scatter matrix:

$$S_W = \sum_{i=1}^C \sum_{m=1}^{N_i} (x_m^i - \mu_i)(x_m^i - \mu_i)^T$$

- Total scatter matrix:

$$S_T = S_B + S_W = \sum_{i=1}^C \sum_{m=1}^{N_i} (x_m^i - \mu)(x_m^i - \mu)^T$$

Road map

- FLD Modifications
- Discriminative Common Vector Method (DCV)
- Rough Common Vector Method (RCV)
- Discriminative Common Vector with Kernels (KDCV)

Outline

- 1 Introduction
 - Background
 - **FLD's modifications**
- 2 Discriminative Common Vector (DCV)
 - Introduction
 - DCV by using the Range Space of S_W
 - DCV by using Difference Subspaces and the Gram-Schmidt Orthogonalization Procedure
- 3 Rough DCV
- 4 Kernel DCV

Some FLD-based methods

- Pseudoinverse method: replacing S_W^{-1} by its pseudoinverse.
- Perturbation method: adding a small perturbation matrix Δ to S_W in order to make it non-singular.
- Rank Decomposition method: making successive eigen-decompositions of the total scatter matrix S_T and the between class scatter matrix S_B .
- They are computationally expensive since the scatter matrices are very large.

Fisherface method

- Two stage method: PCA + Linear Discriminant Analysis.
- PCA is used to reduce data dimensionality so as to make S_W non-singular.
- By PCA use some directions corresponding to the small eigenvalues of S_T are thrown away, removing dimensions with potential discriminative information.

Null Space Method

- Based on the modified FLD criterion:

$$J_{MFLD}(W_{opt}) = \arg \max_W \frac{|W^T S_B W|}{|W^T S_T W|}$$

- This method has been proposed to be used when the dimension of the sample space is larger than the rank of S_W .
- The MFLD criterion attains its maximum when all image samples are projected onto the null space of S_W , and then PCA is applied to the projected samples to obtain the optimal projection vectors.
- The performance of the Null Space method improves if the null space of S_W is large.
- There is not an efficient algorithm for applying this method in the original sample space.

PCA + Null Space Method

- PCA is applied to remove the null space of S_T , which contains the intersection of the null spaces of S_W and S_B .
- Then the optimal projection vectors are found in the remaining lower dimensional space by Null Space method.
- The difference with the Fisherface method is that, here S_W is typically singular in the reduced space because all eigenvectors corresponding to the non-zero eigenvalues of S_T are used for dimension reduction.

Direct-LDA method

- Uses the simultaneous diagonalization method.
- First, the null space of S_B is removed and then, the projection vectors that minimize S_W in the transformed space are selected from the range space of S_B .
- Removing the null space of S_B by dimensionality reduction will also remove part of the null space of S_W removing important discriminant information.
- Furthermore, the whitening process over S_B is redundant.

Comparisons

Table: Comparisons of performance across methods for $n > C - 1$

Rank	Accuracy	Training Time	Testing Time	Storage Requirements
1	DCV, PCA + Null Space	Direct-LDA	DCV, PCA + Null Space	DCV, PCA + Null Space
2	Fisherface	DCV	Fisherface, Direct-LDA	Fisherface, Direct-LDA
3	Direct-LDA	Eigenface	Eigenface	Eigenface
4	Eigenface	Fisherface		
5		PCA + Null Space		

Introduction

- DCV addresses the limitations of previous methods that use the null space of S_W to find the optimal projection vectors.
- It can be only used when the dimension of the sample space is larger than the rank of S_W .
- This approach extracts the common properties of classes in the training set by eliminating the differences of the samples in each class.

Algorithms

- Previous works in word recognition obtain a common vector for each class by removing all the features in the direction of the eigenvectors corresponding to the non-zero eigenvalues of the scatter matrix of its own class.
- Cevikalp's work describes two algorithms to obtain DCV for face recognition:
 - Instead of using a given class's own scatter matrix, he uses the within-classes scatter matrix of all classes to obtain the common vector.
 - He gives an alternative algorithm based on the subspace methods and the Gram-Schmidt orthogonalization procedure.

Outline

- 1 Introduction
 - Background
 - FLD's modifications
- 2 **Discriminative Common Vector (DCV)**
 - Introduction
 - **DCV by using the Range Space of S_W**
 - DCV by using Difference Subspaces and the Gram-Schmidt Orthogonalization Procedure
- 3 Rough DCV
- 4 Kernel DCV

S_W Null Space based criterion

- In the special case where $w^T S_W w = 0$ and $w^T S_B w \neq 0$ for all $w \in \mathbb{R}^d \setminus \{0\}$, the modified FLD criterion attains a maximum.
- A projection vector w satisfying the above conditions does not necessarily maximizes the between-class scatter matrix. In this case, a better criterion is given by:

$$J(W_{opt}) = \arg \max_{|W^T S_W W = 0|} |W^T S_B W| = \arg \max_{|W^T S_W W = 0|} |W^T S_T W|$$

Direct algorithm

- To find the optimal projection vectors w in the null space of S_W , the face samples are projected onto the null space of S_W and then, the projection vectors are obtained by PCA.
- However, this task is computationally intractable since the dimension of the null space can be very large.
- A more efficient way of doing it is by using the orthogonal complement of the null space of S_W , which typically is significantly lower-dimensional space.

Feasible algorithm

Description

- Let R^d be the original sample space, V be the range space of S_W , and V^\perp be the null space of S_W :

$$V = \text{span} \{ \alpha_k | S_W \alpha_k \neq 0, \quad k = 1, \dots, r \}$$

$$V = \text{span} \{ \alpha_k | S_W \alpha_k = 0, \quad k = r + 1, \dots, d \}$$

- Where:
 - $r < d$ is the rank of S_W
 - $\{ \alpha_1, \dots, \alpha_d \}$ is an orthonormal set, and $\{ \alpha_1, \dots, \alpha_r \}$ is the set of orthonormal eigenvectors corresponding to the non-zero eigenvalues of S_W .

Feasible algorithm

Goal

- Considering the matrices $Q = [\alpha_1 \dots \alpha_r]$ and $\tilde{Q} = [\alpha_{r+1} \dots \alpha_d]$.
- Since $R^d = V \oplus V^\perp$, every face image $x_m^i \in R^d$ has a unique decomposition of the form

$$x_m^i = y_m^i + z_m^i$$

- where $y_m^i = Px_m^i = QQ^T x_m^i \in V$, $z_m^i = \tilde{P}x_m^i = \tilde{Q}\tilde{Q}^T x_m^i \in V^\perp$ and P and \tilde{P} are the projection operators onto V and V^\perp respectively.
- The goal is to compute:

$$z_m^i = x_m^i - y_m^i = x_m^i - Px_m^i$$

Common vectors

- The eigenvectors can be obtained from the M by M matrix, $A^T A$ where A is a d by M matrix of the form

$$A = [x_1^1 - \mu_1 \dots x_N^1 - \mu_1 \ x_1^2 - \mu_2 \dots x_N^2 - \mu_2 \dots x_N^C - \mu_C]$$

- Let λ_k and v_k be the k th non-zero eigenvalue and the corresponding eigenvector of $A^T A$. Then, $\alpha_k = Av_k$ will be the eigenvector that corresponds to the k th non-zero eigenvalue of S_W .
- It turns out that we obtain the same unique vector for all samples of the same class, which are defined as the common vectors:

$$x_{com}^i = x_m^i - QQ^T x_m^i = \tilde{Q}\tilde{Q}^T x_m^i, \quad m = 1, \dots, N; i = 1, \dots, C$$

Maximizing criterion

- After obtaining the common vectors x_{com}^i , optimal projection vectors will be those that maximize the scattering of the common vectors:

$$J(W_{opt}) = \arg \max_{|W^T S_W W = 0|} |W^T S_B W| = \arg \max_{|W^T S_W W = 0|} |W^T S_T W| = \arg \max_W |W^T S_T W|$$

- W is a matrix whose columns are the orthonormal optimal projection vectors w_k , and S_{com} is the scatter matrix of the common vectors

$$S_{com} = \sum_{i=1}^C (x_{com}^i - \mu_{com})(x_{com}^i - \mu_{com})^T, \quad i = 1, \dots, C$$

Obtaining optimal projections

- All eigenvectors corresponding to the non-zero eigenvalues of S_{com} will be the optimal projection vectors.
- Instead of using S_{com} that is typically a large d by d matrix, the smaller matrix $A_{com}^T A_{com}$ of size C by C can be used, where

$$A_{com} = [x_{com}^1 - \mu_{com} \dots x_{com}^C - \mu_{com}]$$

- Each class is discriminated by a discriminative common vector:

$$\Omega_i = W^T x_m^i, \quad m = 1, \dots, N; i = 1, \dots, C$$

- To recognize a test image x_{test} , the feature vector of this image is found by $\Omega_{test} = W^T x_{test}$, and the Euclidean distance to each class's discriminative common vector gives the classification.




Algorithm

- Step 1: compute the non-zero eigenvalues and corresponding eigenvectors of S_W by using the matrix $A^T A$. Set $Q = [\alpha_1 \dots \alpha_r]$ where r is the rank of S_W .
- Step 2: choose any sample from each class and project it onto the null space of S_W to obtain the common vectors.
- Step 3: compute the eigenvectors w_k with non-zero eigenvalues of the matrix $A_{com}^T A_{com}$. Use these eigenvectors to form the projection matrix $W = [w_1 \dots w_{C-1}]$.

Outline

- 1 Introduction
 - Background
 - FLD's modifications
- 2 Discriminative Common Vector (DCV)
 - Introduction
 - DCV by using the Range Space of S_W
 - DCV by using Difference Subspaces and the Gram-Schmidt Orthogonalization Procedure
- 3 Rough DCV
- 4 Kernel DCV

For Further Reading I

-  Discriminative Common Vectors for Face Recognition. Hakan Cevikalp, Marian Neamtu, Mitch Wilkes, Atalay Barkana. IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol.27, N^o 1, pp: 4-13. January 2005.
-  Discriminative Common Vector Method with Kernels. Hakan Cevikalp, Marian Neamtu, Mitch Wilkes. IEEE Transactions on Neural Networks, Vol.17, N^o 6, pp: 1550-1565. November 2006.
-  The Kernel Common Vector Method: A Novel Nonlinear Subspace Classifier for Pattern Recognition. Hakan Cevikalp, Marian Neamtu, Atalay Barkana. IEEE Transactions on Systems, Man, and Cybernetics, Vol.37, N^o 4, pp: 937-951. August 2007.

For Further Reading II



Extracción de Características Mediante Vectores Discriminantes
Extendidos con Kernel. Inteligencia Artificial.
<http://erevista.aepia.org>. On printing. 2009.

Questions?

Thank you very much for your attention.

- Contact:
 - Miguel Angel Veganzones
 - Grupo Inteligencia Computacional
 - Universidad del País Vasco - UPV/EHU (Spain)
 - E-mail: miguelangel.veganzones@ehu.es
 - Web page: <http://www.ehu.es/computationalintelligence>