

Multi-view Linear Discriminant Analysis Network (MvLDAN)

**Sawan Patidar (0801CS1717070)
December 13, 2020**

Contents

1. Introduction	2
1.1. CCA	2
1.2. LDA	2
1.3. MvLDAN	2
2. Mathematical Formulation	5
2.1. Formulation	5
2.2. Optimization	8
3. Algorithm	9
3.1. MvLDAN Algorithm	9
4. Documentation of API	10
4.1. Package Organization	10
4.2. Methods	10
5. Learning Outcomes	12
A. References	13

Chapter 1

Introduction

1.1 CCA:

Canonical correlation analysis finds collections linear combinations of pairs of multivariate data vectors that are maximally correlated with each other, and uncorrelated with other pairs. Simultaneously observed data pairs can be analyzed to diagnose relationships between the pairs, and CCA on time-lagged pairs can be used for forecasting. MCA is a related method that maximizes covariance.

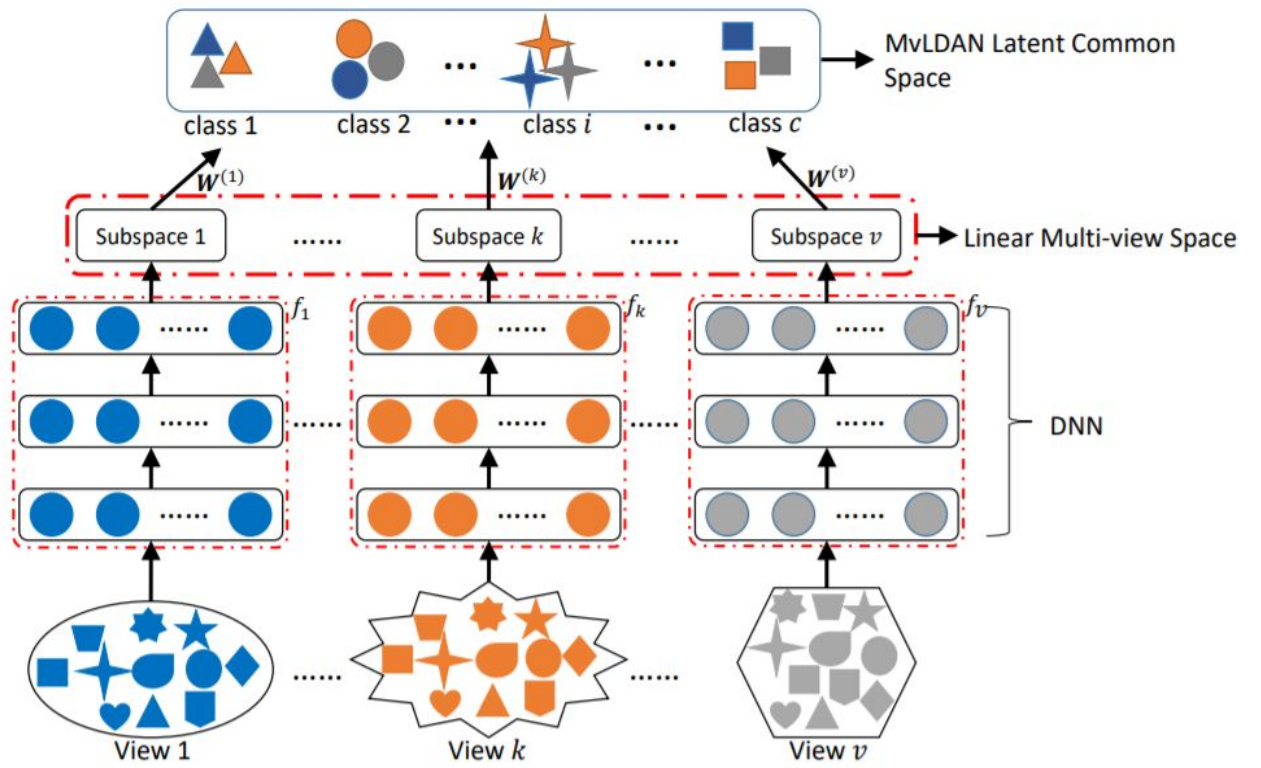
1.2 LDA

Linear Discriminant Analysis or Normal Discriminant Analysis or Discriminant Function Analysis is a dimensionality reduction technique which is commonly used for the supervised classification problems. It is used for modeling differences in groups i.e. separating two or more classes. It is used to project the features in higher dimension space into a lower dimension space.

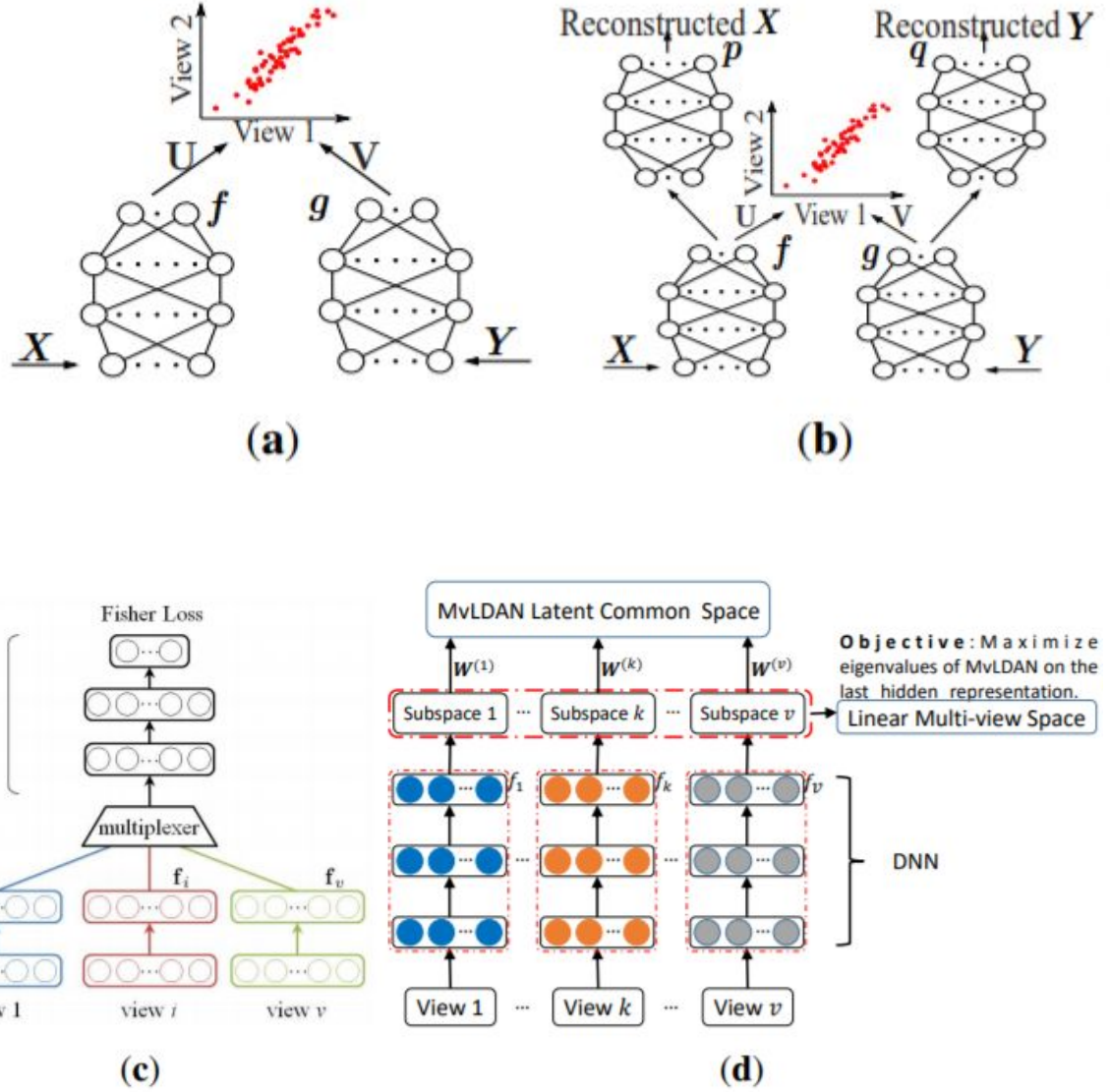
1.3 MvLDAN

To eliminate the complicated (usually highly nonlinear) view discrepancy for favorable cross-view recognition and retrieval, Multi-view Linear Discriminant Analysis Network (MvLDAN) seeks a nonlinear discriminant and view-invariant representation shared among multiple views. Unlike existing multi-view methods which directly learn a common space to reduce the view gap, MvLDAN employs multiple feedforward neural networks (one for each view) and a novel eigenvalue-based multi-view objective function to encapsulate as much discriminative variance as possible into all the available common feature dimensions. With the proposed objective function, the MvLDAN could produce representations possessing:

- a) low variance within the same class regardless of view discrepancy,
- b) high variance between different classes regardless of view discrepancy,
- c) high covariance between any two views.



The above figure shows the Framework of MvLDAN, where different shapes denotes different classes, and different colors represent different views.



The above figure shows the Schematic diagram of DNN-based multi-view learning models. Fig. (a) shows the network model of DCCA. Fig. (b) shows the network model of DCCAE. Fig. (c) shows the network model of MvDN. Fig. (d) shows the network model of MvLDAN. In the figures, X and Y represent different views; f , g , g_c and f_i are the corresponding sub-networks; U , V and $W(k)$ are the view-specific transforms.

Chapter 2

Mathematical Formulation

2.1 Formulation

Let $\mathbf{X}^{(k)} = \{\mathbf{x}_{ij}^{(k)} \in \mathbb{R}^{d_k \times 1} | i = 1, \dots, c; j = 1, \dots, N_i^{(k)}; k = 1, \dots, v\}$ be the samples from the k-th view, where $\mathbf{x}_{ij}^{(k)}$ denotes the j-th sample from the k-th view of the i-th class with d_k dimensionality, c is the number of classes and $N_i^{(k)}$ is the number of samples from the k-th view of the i-th class. As shown in Fig. 1, the sub-network for the k-th view can be represented as a nonlinear function

$f_k(\cdot; \theta_{f_k}) \in \mathbb{R}^{p \times 1}$, where θ_{f_k} is the parameter of the k-th view. Therefore, the output of the k-th view is formulated as

$$\mathbf{y}_{ij}^{(k)} = f_k(\mathbf{x}_{ij}^{(k)}) \quad \text{--- Eq.1}$$

for the sample $\mathbf{x}_{ij}^{(k)}$. The proposed method seeks to find a set of linear transformations $\{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)}, \dots, \mathbf{W}^{(v)}\}$ to project the samples from v views into a latent common space, where $\mathbf{W}^{(k)} \in \mathbb{R}^{p \times q}$ denotes the mapping of the i-th view, and $q = c - 1$ (since \mathbf{S}_b is of rank $c - 1$ at most).

For ease of representation, we denote the projection results by

$$\mathbf{Z} = \{\mathbf{z}_{ij}^{(k)} = \mathbf{W}^{(k)T} \mathbf{y}_{ij}^{(k)} | i = 1, \dots, c; j = 1, \dots, N_i^{(k)}; k = 1, \dots, v\}$$

Unlike the traditional LDA-like methods we use the pairwise-view information to further eliminate the view discrepancy. Mathematically, the objective function is formulated as:

$$\arg \max_{\substack{f_1, f_2, \dots, f_v \\ \mathbf{W}^{(1)}, \dots, \mathbf{W}^{(v)}}} \frac{\text{Tr}(\mathbf{S}_b) + \lambda \text{Tr}(\mathbf{S}_c)}{\text{Tr}(\mathbf{S}_w) + \beta \sum_{k=1}^v \|\mathbf{W}^{(k)}\|_F^2}, \quad \text{--- Eq. 2}$$

where $\text{Tr}(\cdot)$ is the trace operator, \mathbf{S}_b and \mathbf{S}_w correspond to the between-class and within-class scatter matrices from the classes of all views, \mathbf{S}_c is the pairwise-view covariance matrix over all views.

Moreover, the between-class scatter S_b , within-class scatter matrix S_w , and the covariance matrix S_c are defined by:

$$\begin{aligned}
 S_b &= \sum_{i=1}^c N_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T, \\
 S_w &= \sum_{i=1}^c \sum_{k=1}^v \sum_{j=1}^{N_i^{(k)}} (z_{ij}^{(k)} - \boldsymbol{\mu}_i)(z_{ij}^{(k)} - \boldsymbol{\mu}_i)^T, \\
 S_c &= \sum_{k < l}^v \mathbf{Z}^{(k)} \mathbf{Z}^{(l)T},
 \end{aligned}
 \quad \text{--- Eq 3,4,5}$$

Where $\boldsymbol{\mu}_i = \frac{1}{N_i^{(k)}} \sum_{k=1}^v \sum_{j=1}^{N_i^{(k)}} z_{ij}^{(k)}$ is the mean of all samples in the i-th class from the k-th view, $\boldsymbol{\mu} = \frac{1}{N} \sum_{k=1}^v \sum_{i=1}^c \sum_{j=1}^{N_i^{(k)}} z_{ij}^{(k)}$ is the mean of all samples from all views.

The between-class scatter matrix in the common space can be further written as

$$S_b = \mathbf{W}^T \begin{bmatrix} B_{11} & \cdots & B_{1v} \\ \vdots & \ddots & \vdots \\ B_{v1} & \cdots & B_{vv} \end{bmatrix} \mathbf{W} = \mathbf{W}^T \mathbf{B} \mathbf{W}, \quad \text{--- Eq. 6}$$

where $\mathbf{W} = [\mathbf{W}^{(1)T} \ \mathbf{W}^{(2)T} \ \cdots \ \mathbf{W}^{(v)T}]^T$ is constructed from the transformation matrices for all views and B_{kl} is defined as:

$$B_{kl} = \sum_{i=1}^c \frac{1}{N_i} \mathbf{s}_i^{(k)} \mathbf{s}_i^{(l)T} - \frac{1}{N} \mathbf{s}^{(k)} \mathbf{s}^{(l)T}, \quad \text{--- Eq. 7}$$

Where $\mathbf{s}_i^{(k)} = \sum_{j=1}^{N_i^{(k)}} \mathbf{y}_{ij}^{(k)}$ is the sum of all samples in the i-th class from the k-th view; $\mathbf{s}^{(k)} = \sum_{i=1}^c \mathbf{s}_i^{(k)}$ is the sum of all samples from the k-th view.

Similarly other equations can be re written as $S_w = \mathbf{W}^T \mathbf{C} \mathbf{W}$ and $S_c = \mathbf{W}^T \mathbf{D} \mathbf{W}$, where \mathbf{C} and \mathbf{D} are partitioned matrices like \mathbf{B} , the (k,l)th submatrices of \mathbf{C} and \mathbf{D} and are defined as:

$$\begin{aligned}
 C_{kl} &= \sum_{i=1}^c \left((k == l) \sum_{j=1}^{N_i^{(k)}} \mathbf{y}_{ij}^{(k)} \mathbf{y}_{ij}^{(l)T} - \frac{1}{N_i} \mathbf{s}_i^{(k)} \mathbf{s}_i^{(l)T} \right) \\
 D_{kl} &= \begin{cases} \mathbf{0} & k = l \\ \mathbf{Y}^{(k)} \mathbf{Y}^{(l)T} & k \neq l \end{cases},
 \end{aligned}
 \quad \text{--- Eq. 8,9}$$

where $k == l$ is a Boolean equation, whose value is 1 if $k = l$ and 0 otherwise; $\mathbf{0}$ is a zero matrix. Thus, our objective function can be rewritten as

$$\arg \max_{\substack{f_1, f_2, \dots, f_v \\ \mathbf{W}^{(1)}, \dots, \mathbf{W}^{(v)}}} \frac{\text{Tr}(\mathbf{W}^T (\mathbf{B} + \lambda \mathbf{D}) \mathbf{W})}{\text{Tr}(\mathbf{W}^T (\mathbf{C} + \beta \mathbf{I}) \mathbf{W})}, \quad \text{--- Eq. 10}$$

where \mathbf{I} denotes the identity matrix.

The equation can be relaxed as the following tractable ratio trace (determinant ratio) form:

$$\arg \max_{\substack{f_1, f_2, \dots, f_v \\ \mathbf{W}^{(1)}, \dots, \mathbf{W}^{(v)}}} \text{Tr}((\mathbf{S}_w + \beta \mathbf{W}^T \mathbf{W})^{-1} (\mathbf{S}_b + \lambda \mathbf{S}_c)) = \arg \max_{\substack{f_1, f_2, \dots, f_v \\ \mathbf{W}^{(1)}, \dots, \mathbf{W}^{(v)}}} \frac{|\mathbf{W}^T (\mathbf{B} + \lambda \mathbf{D}) \mathbf{W}|}{|\mathbf{W}^T (\mathbf{C} + \beta \mathbf{I}) \mathbf{W}|}, \quad \text{--- Eq. 11}$$

where $|\cdot|$ is the determinant operator. The above eq is s equivalent to the following generalized eigenvalue decomposition (GED) problem

$$(\mathbf{B} + \lambda \mathbf{D}) \mathbf{w}_i = \gamma_i (\mathbf{C} + \beta \mathbf{I}) \mathbf{w}_i, \quad \text{--- Eq. 12}$$

Where $\gamma_i |_{i=1}^{c-1}$ is the i -th largest eigenvalue of the GED with the corresponding eigenvector \mathbf{w}_i , and \mathbf{w}_i constitutes the i -th column vector of the matrix \mathbf{W} .

Moreover $\gamma_i |_{i=1}^{c-1}$ and $\mathbf{w}_i |_{i=1}^{c-1}$ are the eigenvalues and eigenvectors of $(\mathbf{C} + \beta \mathbf{I})^{-1} (\mathbf{B} + \lambda \mathbf{D})$, respectively.

Theorem 1: : Let γ_i be the i -th largest eigenvalue of the GED Eq given below, then the DNN's $\mathbf{f}^* = \{f_1^*, f_2^*, \dots, f_v^*\}$ learnt by MvLDAN will be the solution to Eq 11.

$$\arg \max_{f_1, f_2, \dots, f_v} \text{Tr}((\mathbf{S}_w + \beta \mathbf{W}^T \mathbf{W})^{-1} (\mathbf{S}_b + \lambda \mathbf{S}_c)). \quad \text{---Eq. 13}$$

Equivalently, \mathbf{f}^* is also solution to

$$\arg \max_{f_1, f_2, \dots, f_v} \sum_{i=1}^{c-1} \gamma_i. \quad \text{--- Eq. 14}$$

A certain threshold is used to filter the relatively larger eigenvalues when optimizing the parameters of the DNNs, i.e. we just consider the m eigenvalues that do not exceed the threshold for variance maximization:

$$\begin{aligned} \arg \min_{f_1, f_2, \dots, f_v} & -\frac{1}{m} \sum_{i=1}^m \gamma_i \quad \text{with} \\ \{\gamma_1, \dots, \gamma_m\} &= \{\gamma_i | \gamma_i < \min\{\gamma_1, \dots, \gamma_{c-1}\} + \epsilon\}. \end{aligned} \quad \text{--- Eq. 19}$$

2.2 Optimisation

1) Feedforward and calculate the loss: For the k -th view, the samples of a batch of $X^{(k)}$ are fed forward to the MvLDAN as in Eq. (1), and the output of the MvLDAN is denoted as $Y^{(k)}$, with $Y^{(k)} = \{y_{ij}^{(k)} | i = 1, \dots, c; j = 1, \dots, N_b\}$. Where N_b is the batch size. The loss of the whole network is calculated as in Eq. (19), denoted as $\mathcal{J} = -\frac{1}{m} \sum_{i=1}^m \gamma_i$.

2) Gradient of the loss layer: As the loss J is directly calculated with the outputs of the network, there is no parameter involved. We firstly need to calculate the gradient of \mathcal{J} with respect to $Y^{(k)}$. If $w_i^T (C + \beta I) w_i = 1$, then,

$$\frac{\partial \gamma_i}{\partial Y^{(k)}} = w_i^T \left(\frac{\partial (B + \lambda D)}{\partial Y^{(k)}} - \gamma_i \frac{\partial (C + \beta I)}{\partial Y^{(k)}} \right) w_i. \quad \text{--- Eq. 20}$$

Therefore the gradient of \mathcal{J} w.r.t. $Y^{(k)}$ can be calculated as follows

$$\frac{\partial \mathcal{J}}{\partial Y^{(k)}} = \frac{1}{m} \sum_{i=1}^m w_i^T \left(\frac{\partial (B + \lambda D)}{\partial Y^{(k)}} - \gamma_i \frac{\partial (C + \beta I)}{\partial Y^{(k)}} \right) w_i. \quad \text{--- Eq. 21}$$

3. MvLDAN update via gradient descent: Let us denote the parameters of the k -th sub-network as θ_{f_k} , then θ_{f_k} is updated by descending the stochastic gradient as follows:

$$\nabla_{\theta_{f_k}} \mathcal{J} = \frac{\partial \mathcal{J}}{\partial Y^{(k)}} \cdot \frac{\partial Y^{(k)}}{\partial f_k}. \quad \text{--- Eq. 22}$$

Chapter 3

Algorithm

This section we will come to know about the algorithm used in Multi-view Linear Discriminant Analysis Network.

3.1 MvLDAN Algorithm:

Algorithm 1 Optimization procedure of MvLDAN

Input: Training set $\mathcal{X} = \{\mathbf{X}^{(k)}\}_{k=1}^v$, the corresponding class label set, batch size N_b , and learning rate α .

- 1: **while** not converge **do**
- 2: Randomly select N_b samples for each view from \mathcal{X} to construct a multi-view mini-batch.
- 3: Compute the outputs of each view-specific sub-network on the multi-view mini-batch according to Eq. (1) as $\mathbf{y}_{ij}^{(k)} = f_k(\mathbf{x}_{ij}^{(k)})$.
- 4: Compute \mathbf{B} , \mathbf{D} and \mathbf{C} on the obtained outputs according to Eq. (7), Eq. (8) and Eq. (9), respectively.
- 5: Compute the eigenvalues according to Eq. (12).
- 6: Update the parameters of each sub-network by minimizing \mathcal{J} in Eq. (19) by descending their stochastic gradient according to Eq. (22) as $\theta_{f_k} = \theta_{f_k} - \alpha \nabla_{\theta_{f_k}} \mathcal{J}$.
- 7: **end while**
- 8: Compute the transformations $\{\mathbf{W}^{(k)}\}_{k=1}^v$ on \mathcal{X} according to Eq. (1), Eq. (7), Eq. (8), Eq. (9) and Eq. (12).

Output: Optimized MvLDAN model.

Chapter 4

Documentation of API

4.1 Package Organisation:

```
from MvLDAN import MvLDAN_general
```

```
loss_out = Lambda(MvLDAN_gneral, output_shape=(1,),  
name='ctc')(net_output + net_labels)
```

4.2 Methods:

load_noisyMNIST():

Load the dataset in training, validating and testing sets.

create_mnist_full_model((input_size, output_size, value_l2, learning_rate)

Create a model from the mnist dataset:

Parameters:

input_size = ndarray(n*p)

Where p is feature and n is number of samples

output_size = ndarray(n*p)

Where p is feature and n is number of samples

Value_l2 = constant with value 1e-5

learning_rate = constant with value 1e-3

th_MvLDAN_Sw_Sb(data, labels):

Method to find out the within and between scatter matrix.

Parameters:

data = training data (50000 inputs)

labels = labels associated with training data

th_MvLDAN(data_inputs, labels):

Method returns the mean, standard deviation, the proposed set of linear transformations and the eigen values.

Parameters:

data = training data (50000 inputs)

labels = labels associated with training data

Chapter 5

Learning Outcomes

- Learned about CCA and LDA
- Feedforward Neural Networks.
- Limitations of various unsupervised and supervised approaches and advantages of MvLDAN over them.
- Eliminating various complicated view discrepancy for favourable cross view recognition and retrieval.
- How to study a research paper

Chapter 6

References

- Deep Learning: Feedforward Neural Network
<https://towardsdatascience.com/deep-learning-feedforward-neural-network-26a6705dbdc7>
- Multi-view Linear Discriminant Analysis Network
<https://ieeexplore.ieee.org/document/8704986>