

End-term Project Report : Joint feature selection and subspace learning

*Team Name: Vikings**Team Members: 200260040,200260045***Abstract**

Dealing with data with large feature space comes along with the curse of dimensionality. Hence, there is a need for dimensionality reduction. Dimensionality reduction can be done by feature selection or subspace learning. In this project, we attempt to understand and investigate a novel framework for subspace learning which attempts to perform subspace learning and feature selection simultaneously by using a sparsity inducing objective function.

1 Introduction

Dimensionality reduction algorithms can be roughly classified into two categories: Feature selection and subspace learning. Feature selection involves selecting the relevant features and eliminating the rest while in subspace learning we attempt to learn a subspace by taking combinations of features. Examples of well-known subspace learning algorithms are Principal Component analysis, Linear Discriminant Analysis etc. It has been shown that such algorithms can be incorporated into a common framework by using a graph embedding approach. There also exist variations of the above algorithms which can be used to learn sparse sub spaces. However in most of these cases, the selected features where different for different samples (depending on the pattern of sparsity). In this framework we attempt to learn sub spaces which are sparse along features i.e have feature value as 0 for all samples using the properties of the $l_{2,1}$ norm.

In the remaining part of this project, we look into relevant literature before briefly introducing the graph embedding framework and looking at an objective function for joint feature selection and subspace learning (FSSL) based on it. We then give details of the experiments followed by concluding remarks.

2 Literature Survey

Our project is based on the algorithm proposed by the paper [2]. It draws inspiration from subspace learning algorithms proposed by PCA paper and LDA paper as well as sparse subspace learning algorithms [4]. The proposed algorithm is based on a graph embedding framework proposed by [3]. In this paper, a general view for subspace learning has been provided using the concept of graph embeddings wherein we construct a graph using data points connected with each other by a weights represented as an adjacency matrix. They have shown equivalence of this view with well known linear subspace learning algorithms like PCA and LDA. Convergence guarantees of the proposed algorithms and relevant mathematical background has been given in [1].

3 Methods and Approaches

3.1 Work done before mid-term project review

3.1.1 Notations

We denote the data matrix $\mathbf{X}=[\mathbf{x}_1, \mathbf{x}_2 \dots \mathbf{x}_n] \in \mathbb{R}^{d \times n}$. The aim is to learn a projection matrix $\mathbf{A} \in \mathbb{R}^{d \times m}$ thus projecting the data to an m-dimensional subspace. The i th row of \mathbf{A} is denoted by \mathbf{a}^i and the j th column of \mathbf{A} is denoted by \mathbf{a}_j . The Frobenius norm of \mathbf{A} is defined as $\|\mathbf{A}\|_F = \sqrt{\sum_i \|\mathbf{a}^i\|_2^2}$, the $l_{2,1}$ norm is defined as $\|\mathbf{A}\|_{2,1} = \sum_i \|\mathbf{a}^i\|_2$ and the $l_{1,1}$ norm is defined as $\|\mathbf{A}\|_{1,1} = \sum_i \|\mathbf{a}^i\|_1$

3.1.2 Graph embedding view of subspace learning

In the graph embedding view, for linear dimensionality reduction the optimal \mathbf{A} is given by the following optimisation problem:

$$\begin{aligned} \min_{\mathbf{A}} \text{tr}(\mathbf{A}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{A}) \\ \text{where } \mathbf{A}^T \mathbf{X} \mathbf{D} \mathbf{X}^T \mathbf{A} = \mathbb{I} \end{aligned} \quad (1)$$

To do feature selection it is desirable to have some rows of the projection matrix be all zeros. Here $D_{ii} = \sum_j W_{ij}$ is a diagonal matrix and $\mathbf{L}=\mathbf{D}-\mathbf{W}$ is called the graph Laplacian. This problem is solved using generalized eigen-decomposition $\mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{A} = \Lambda \mathbf{X} \mathbf{D} \mathbf{X}^T \mathbf{A}$ where Λ is a diagonal matrix of eigen values. For the linear Graph embedding framework the different choices for \mathbf{W} lead to different linear dimensionality methods like PCA, LDA, LPP etc. For LDA, if there are c classes and the k th class has n_k samples such that $n_1 + n_2 + \dots n_k = n$, then

$$W_{ij} = \begin{cases} \frac{1}{n_k} & \text{if } \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ belong to } k\text{th class} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

For LPP the weight matrix is:

$$W_{ij} = \begin{cases} d(x_i, x_j), & \text{if } \mathbf{x}_j \in \mathcal{N}_k(\mathbf{x}_i) \text{ or } \mathbf{x}_i \in \mathcal{N}_k(\mathbf{x}_j) \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

Here, $\mathcal{N}_k(x_i)$ is the set of k nearest neighbors of \mathbf{x}_i and $d(\mathbf{x}_i, \mathbf{x}_j)$ is a measure of the similarity between the 2 vectors \mathbf{x}_i and \mathbf{x}_j .

3.1.3 Joint Feature selection and Subspace learning

We want row sparsity in the projection matrix. This motivates us to minimise the $l_{2,1}$ norm which promotes row sparsity. Hence, the formulation becomes:

$$\min_{\mathbf{A}} \|\mathbf{A}\|_{2,1} + \text{tr}(\mathbf{A}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{A}) \quad (5)$$

where the constraint $\mathbf{A}^T \mathbf{X} \mathbf{D} \mathbf{X}^T \mathbf{A} = \mathbb{I}$ holds as before. Since, the constraint is not convex, optimising directly is difficult. Hence, we look for an equivalent formulation. For this formulation, the following theorem is useful:

Theorem 3.1 Let $\mathbf{Y} \in \mathbb{R}^{n \times m}$ be a matrix of which each column is an eigenvector of the eigen-problem $\mathbf{W}\mathbf{y} = \lambda\mathbf{D}\mathbf{y}$. If there exists a matrix $\mathbf{A} \in \mathbb{R}^{d \times m}$ such that $\mathbf{X}^T\mathbf{A} = \mathbf{Y}$, then each column of \mathbf{A} is an eigenvector of eigen-problem $\mathbf{X}\mathbf{W}\mathbf{X}^T\mathbf{a} = \lambda\mathbf{X}\mathbf{D}\mathbf{X}^T\mathbf{a}$.

A proof of this theorem is given in (cite). Using this theorem, we can solve the eigen-problem $\mathbf{X}\mathbf{W}\mathbf{X}^T\mathbf{A} = \lambda\mathbf{X}\mathbf{D}\mathbf{X}^T\mathbf{A}$ to get \mathbf{A} as follows:

- Solve the eigen problem $\mathbf{W}\mathbf{Y} = \lambda\mathbf{D}\mathbf{Y}$ to obtain \mathbf{Y}
- Find \mathbf{A} which satisfies $\mathbf{X}^T\mathbf{A} = \mathbf{Y}$

Since $\mathbf{X}^T\mathbf{A} = \mathbf{Y}$ is a linear system, it can have infinitely many solutions, one solution or no solution. Hence depending on this, we can formulate the problem as shown below.

3.1.4 Algorithm 1

When the linear system has infinitely many solutions the proposed formulation is:

$$\min_{\mathbf{A}} \|\mathbf{A}\|_{2,1} \quad (6)$$

where $\mathbf{X}^T\mathbf{A} = \mathbf{Y}$ The lagrangian of the above objective is:

$$L(\mathbf{A}) = \|\mathbf{A}\|_{2,1} - \text{tr}(\mathbf{\Gamma}^T(\mathbf{X}^T\mathbf{A} - \mathbf{Y})) \quad (7)$$

Minimising the lagrangian by setting the derivative with respect to \mathbf{A} to zero implies:

$$\frac{\partial L(\mathbf{A})}{\partial \mathbf{A}} = \mathbf{G}\mathbf{A} - \mathbf{X}\mathbf{\Gamma} = 0 \quad (8)$$

Here, \mathbf{G} is a diagonal matrix where:

$$g_{ii} = \begin{cases} 0 & \text{if } \mathbf{a}^i = \mathbf{0} \\ \frac{1}{\|\mathbf{a}^i\|_2} & \text{otherwise} \end{cases} \quad (9)$$

Assuming \mathbf{G} is invertible, if we multiply both sides of equation (8) by $\mathbf{X}^T\mathbf{G}^{-1}$ and using the constraint $\mathbf{X}^T\mathbf{A} = \mathbf{Y}$ we get:

$$\mathbf{\Gamma} = (\mathbf{X}^T\mathbf{G}^{-1}\mathbf{X})^{-1}\mathbf{Y} \quad (10)$$

If \mathbf{G} is not invertible, we can add a smoothing term ϵ to g_{ii} . If a decreasing value of ϵ is used the algorithm converges to a global solution. Hence we get (by substituting equation 10 into 8)

$$\mathbf{G}\mathbf{A} - \mathbf{X}(\mathbf{X}^T\mathbf{G}^{-1}\mathbf{X})^{-1}\mathbf{Y} = 0 \quad (11)$$

$$\Rightarrow \mathbf{A} = \mathbf{G}^{-1}\mathbf{X}(\mathbf{X}^T\mathbf{G}^{-1}\mathbf{X})^{-1}\mathbf{Y} \quad (12)$$

In the code, we initialise $\mathbf{G} = \mathbb{I}$, then we compute \mathbf{A} based on equation (12) and then using this value we update the value of \mathbf{G} using equation (9). The updated value of \mathbf{G} is then used to update \mathbf{A} (eqn.12) and this process is repeated until convergence. The convergence of this algorithm was proved in (nie et al).

Situation 1 and its code was covered in pre-midsem review

3.2 Work done after mid-term project review

After the mid-term review, we looked into the case when linear system has a unique solution or no solutions.

3.2.1 Algorithm 2

When the linear system has a single unique solution or no solution then also it is a solution of (6). However, it does not have row sparsity. Hence the reformulated problem is:

$$\min_{\mathbf{A}} \|\mathbf{A}\|_{2,1} \quad (13)$$

such that $\|\mathbf{X}^T \mathbf{A} - \mathbf{Y}\|_F^2 \leq \delta$ or equivalently:

$$\min_{\mathbf{A}} \|\mathbf{A}\|_{2,1} + \mu \|\mathbf{X}^T \mathbf{A} - \mathbf{Y}\|_F^2 \quad (14)$$

As before if we set the derivative to zero we get:

$$\mathbf{A} = 2\mu(\mathbf{G} + 2\mu\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{Y} \quad (15)$$

We can simplify this further using the Woodbury Matrix identity:

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1} \quad (16)$$

from which we can obtain:

$$\mathbf{A} = \mathbf{G}^{-1}\mathbf{X}(\mathbf{X}^T\mathbf{G}^{-1}\mathbf{X} + \frac{1}{\mu}\mathbb{I})^{-1}\mathbf{Y} \quad (17)$$

3.2.2 Using other norms for the objective function

Instead of using the $l_{2,1}$ norm we can use other norms like $l_{1,1}$ norm or frobenius norm. We investigate the effects of using these norms in the ipython notebook submission. In such cases, for both algorithm 1 and 2 the only change appears in the form of \mathbf{G} (while differentiating the norm). For Frobenius norm the form of \mathbf{G} becomes:

$$g_{ii} = \begin{cases} 0 & \text{if } \|\mathbf{A}\|_F = 0 \\ \frac{1}{\|\mathbf{A}\|_F} & \text{otherwise} \end{cases} \quad (18)$$

For the l_{11} norm the form of \mathbf{G} becomes:

$$g_{ij} = 1 \text{ for all } i \text{ and } j \quad (19)$$

Note that for l_{11} norm \mathbf{G} is not a diagonal matrix.

4 Data set Details

For the experiments the following dataset were used:

- Yale Face subset: A set of 165 images of 15 subjects. There are 11 images per subject, one for each of the following facial expressions or configurations: 'happy', 'normal', 'surprised', 'sad', 'sleepy', 'wink', 'centerlight', 'leftlight', 'rightlight', 'glasses', 'noglasses'. The data is resized to 32x32 grey images. Each face is represented as a 1024-dimensional vector.
- MNIST-784 dataset: The MNIST database of handwritten digits with 784 features (from 0 to 9). It is a subset of a larger set available from NIST. The digits have been size-normalized and centered in a fixed-size image with dimensions 28x28. Each image is represented as a 784-dimensional vector. For the purposes of our code we test on data from 2 classes '1' and '9'.

5 Experiments

We investigated the FSSL algorithm using both LDA and LPP type adjacency matrices for both the datasets. We compare the performance of FSSL(LDA) with other subspace learning algorithms like PCA and LDA (using scikit-learn) using accuracy. We selected p number of samples from each class as the training set and the rest of the samples as the test set. For dataset 1, $p=5,6,7$ and for dataset 2, $p=10,20,30$. The training set is randomly chosen and hence the experiment is repeated 20 times to calculate average accuracy for FSSL-LDA. The training set was used to learn a subspace after which recognition was done by a 1-Nearest Neighbor classifier. The baseline model is 1-Nearest Neighbor. μ can be selected by tuning and the smoothing term ϵ is set to 0.01. For PCA, the number of components was set to capture 99% variance. We also studied the effect of value of μ on accuracy. We also look into the case where l_{11} and Frobenius norm were used and their effect on accuracy. All tests were done in Python on a 11th Gen Intel(R) Core(TM) i7-11370H @ 3.30GHz, Windows 11 Machine.

6 Results

- For Yale Face database we obtained low accuracy on baseline as well as PCA, LDA and FSSL-LDA. The accuracy was $\approx 16\%$. Possible causes were due to lower samples per class in the data (1 sample per class).
- For MNIST dataset, we obtained high accuracy of $\approx 98\%$ in baseline, as well as PCA, LDA and FSSL-LDA. FSSL-LDA gave competitive results in spite of having the subspace with lowest dimensionality.
- In both cases, it was observed that the value of μ did not have significant correlation with accuracy.
- Using Frobenius norm or $l_{1,1}$ norm in place of $l_{2,1}$ did not have significant differences in accuracy in both cases.

7 Future Work

The above frameworks can be explored in kernel space. We would also like to explore the frameworks on other forms of large data.

8 Conclusion

In summary we looked at a novel algorithm which attempts to perform joint feature selection and subspace learning. We saw the motivation behind this was to maintain the physical importance of features and also saw why we needed to induce row sparsity in the projection matrix. We also saw the significance of choosing $l_{2,1}$ norm and what are the implications when we replaced it with $l_{1,1}$ and Frobenius norm. The resulting framework gave competitive results to well known subspace algorithms like PCA and LDA for a lower dimensional subspace. We also saw that the value of μ did not significantly affect accuracy.

Baseline (KNN, k=1)				Using LDA learnt subspace			
	5	6	7		5	6	7
Accuracy	0.168636	0.173737	0.169318	Accuracy	0.600909	0.598485	0.592045
Time taken	7.208 ms	3.803 ms	4.897 ms	Time taken	3.948 ms	3.759 ms	3.408 ms
Using LDA Classifier				FSSL-LDA			
	5	6	7	value of mu= 0.1			
Accuracy	0.793636	0.792929	0.784091		5	6	7
Time taken	26.359 ms	30.547 ms	28.057 ms	Accuracy	0.170455	0.173737	0.175568
				Time taken	4.091 ms	5.284 ms	3.623 ms
Using PCA				value of mu= 0.01			
	5	6	7		5	6	7
Accuracy	0.168636	0.173737	0.169318	Accuracy	0.172727	0.173737	0.177273
Time taken	3.354 ms	2.884 ms	3.472 ms	Time taken	4.906 ms	4.002 ms	4.138 ms
				value of mu= 0.001			
	5	6	7		5	6	7
Accuracy	0.168636	0.173737	0.169318	Accuracy	0.179545	0.175758	0.167614
Time taken	3.354 ms	2.884 ms	3.472 ms	Time taken	4.234 ms	4.291 ms	3.722 ms

Figure 1: Figure describing accuracies across PCA, LDA, baseline and FSSL-LDA for Yale Face subset

References

- [1] Xiaofei He and Partha Niyogi. Locality preserving projections. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems*, volume 16. MIT Press, 2003.
- [2] Zhenhui Li Quanquan Gu and Jiawei Han. Joint feature selection and subspace learning. In *Twenty-Second International Joint Conference on Artificial Intelligence*.
- [3] Shuicheng Yan, Dong Xu, Benyu Zhang, Hong-jiang Zhang, Qiang Yang, and Stephen Lin. Graph embedding and extensions: A general framework for dimensionality reduction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(1):40–51, 2007.
- [4] Hui Zou, Trevor Hastie, and Robert Tibshirani. Sparse principal component analysis. *Policy*, pages 1–30, 05 2004.

Baseline (KNN, k=1)

	10	20	30
Accuracy	0.967013	0.976725	0.9818
Time taken	824.212 ms	820.323 ms	864.609 ms

PCA

	10	20	30
Accuracy	0.967455	0.977124	0.982098
Time taken	719.809 ms	736.033 ms	720.307 ms

FSSL (LDA)

	10	20	30
Accuracy	0.97001	0.980696	0.985391
Time taken	250.23 ms	257.454 ms	261.497 ms

Figure 2: Figure describing accuracies across PCA, baseline and FSSL-LDA for MNIST dataset

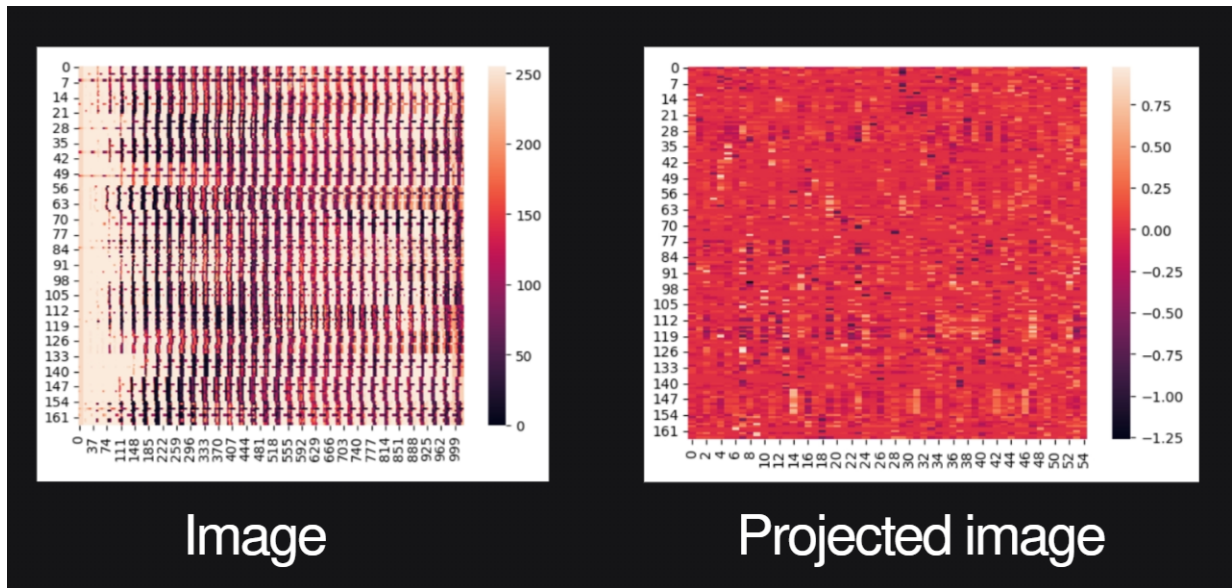


Figure 3: Image and Projected image for Yale Face subset

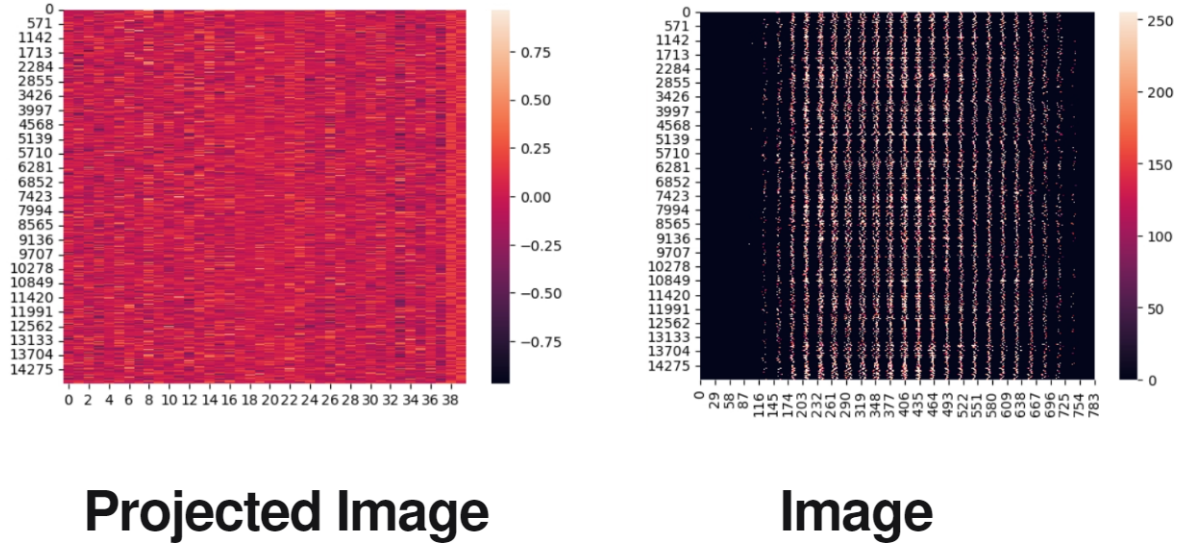
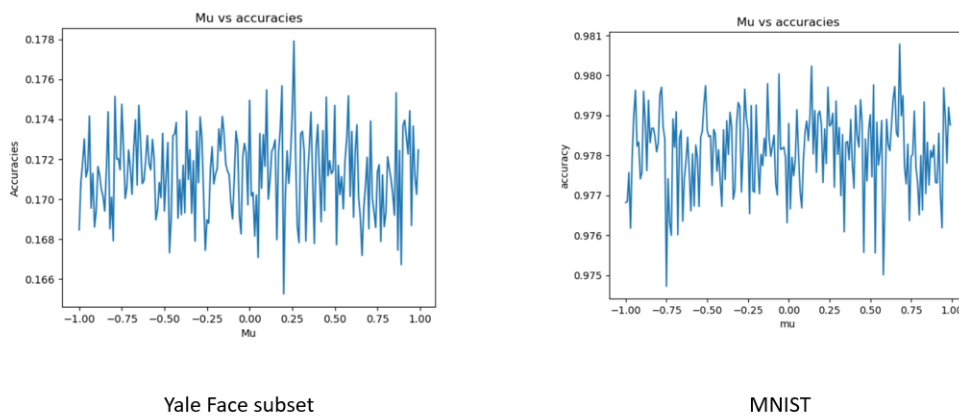


Figure 4: Image and Projected image for MNIST dataset

Figure 5: Plots of μ vs accuracies