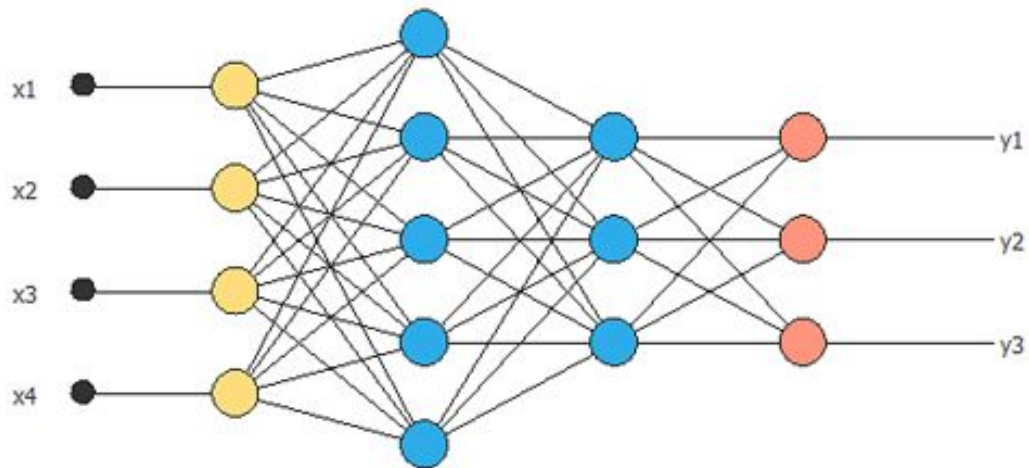


# SMAI ASSIGNMENT-2 REPORT



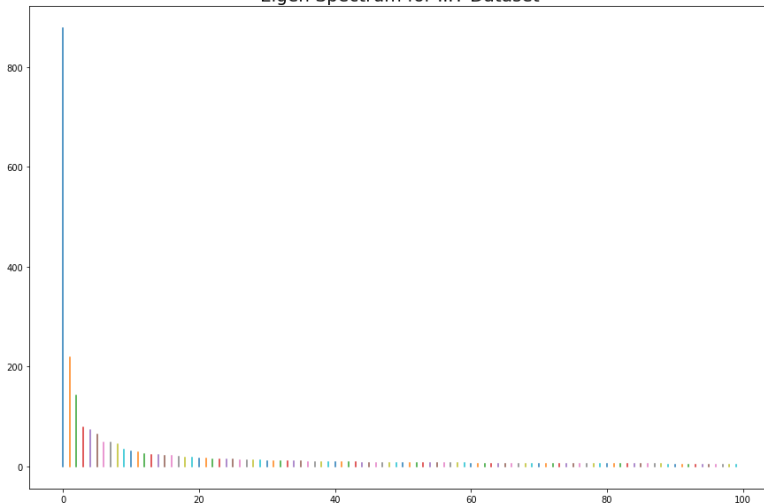
**PURU GUPTA**

**20171187**

## Questions

1. How many eigenvectors/faces are required to “satisfactorily” reconstruct a person in these three datasets?

Eigen Spectrum for IIIT Dataset



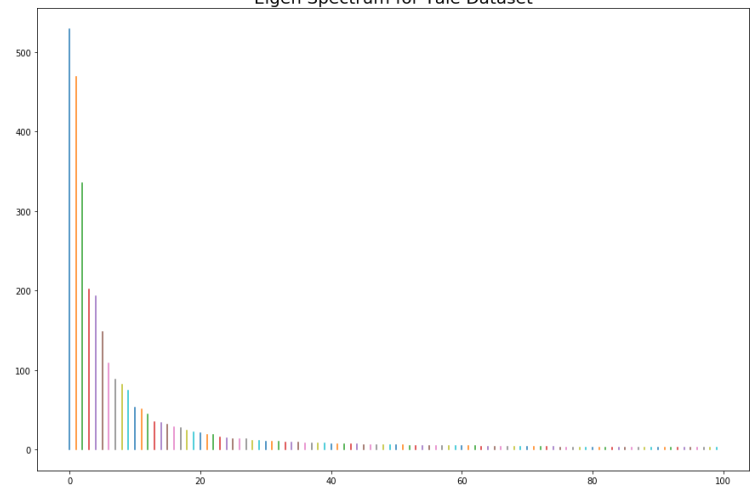
### IIIT Dataset

As seen from the eigen spectrum, 98 eigenvectors are enough to reconstruct a cartoon in IIIT Dataset as it covers around 95% of variance of the dataset and 90% of the sum of total eigenvalues.

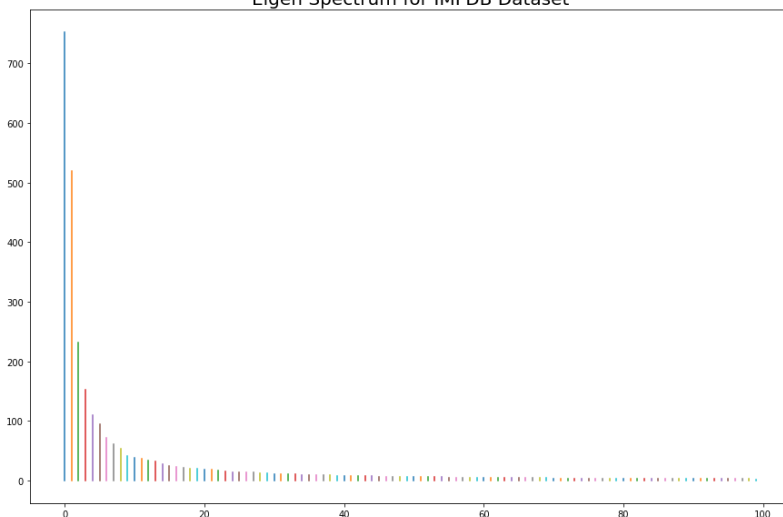
### Yale Dataset

As seen from the eigen spectrum, 13 eigenvectors are enough to reconstruct a person in Yale Dataset as it covers around 95% of variance of the dataset and 90% of the sum of total eigenvalues.

Eigen Spectrum for Yale Dataset



Eigen Spectrum for IMFDB Dataset



### IMFDB Dataset

As seen from the eigen spectrum, 62 eigenvectors are enough to reconstruct a person in IMFDB Dataset as it covers around 95% of variance of the dataset.

## 2. What are eigenfaces?

Representation of all the face images in a dataset as linear combination of the most prominent eigenvectors (higher eigenvalue signifies more prominence) of the covariance matrix of the data matrix is known as eigenface technique.

We try to represent all the data samples with the help of a few eigenfaces.

Example :- Sample face X would be 10% Eigenface A, 13.8 % Eigenface B, and so on, while someone else's face would have a different combination of those same eigenfaces.

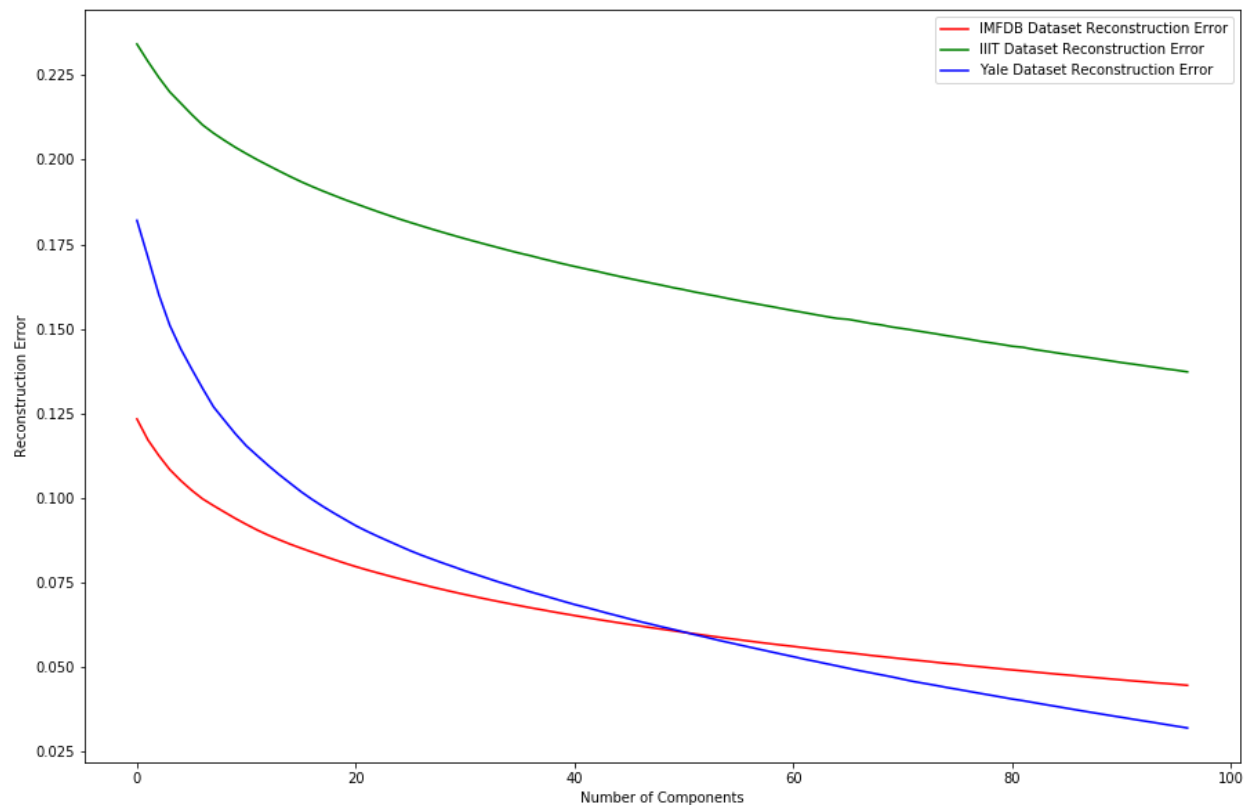
## 3. Which person/identity is difficult to represent compactly with fewer eigenvectors? Why is that? Explain with your empirical observations and intuitive answers.

We reconstructed the images in IMFDB Dataset from 10 principal components and calculated the mean reprojection error for all the 8 labels and it is found that images corresponding to label = 7 (Amir) have the highest mean error, so it is difficult to represent Amir's pictures compactly with fewer eigenvectors.

Reconstruction Error for different classes in IMFDB Dataset

index	error
0.0	0.0051
1.0	0.0086
2.0	0.0083
3.0	0.0115
4.0	0.0117
5.0	0.0105
6.0	0.009
7.0	0.0117

## 4. Which dataset is difficult to represent compactly with fewer eigenvectors? Why is it so? Explain with your empirical observations and intuitive answers.



It can be seen from the above plot that the reconstruction error for the IIIT-CFW class is the highest for any number of features, due to the fact that cartoons don't have predefined facial features and have unique features per sample.

## 5. Which method works well for classification?

To classify the samples in each dataset, a total of 15 combinations of feature extraction and classifiers are used, out of which the top 6 combinations which predicted with the highest accuracy have been reported in the tables below.

Train Set = 80% of the samples

Test Set = 20% of the samples

The feature extraction methods and classifiers are reported below :-

Feature extraction methods	Classifiers
PCA	SVM, Logistic Regression, MLP

LDA	SVM, Logistic Regression, MLP
ResNet	SVM, Logistic Regression, MLP
Resnet + KPCA	SVM, Logistic Regression, MLP
VGG + PCA	SVM, Logistic Regression, MLP

## IMFDB Dataset

For IMFDB Dataset, the combination of LDA with all the three classifiers worked best, giving an accuracy of 100% on the test set.

Classifier Method vs Accuracy Table for IMFDB Dataset

Method	Reduced Space	Classification Error	Accuracy	F1 Score
LDA + SVM	30	0.0	1.0	1.0
LDA + LR	30	0.0	1.0	1.0
ResNet + SVM	2048	0.570087712549569	0.975	0.9770792932557639
(ResNet + KPCA) + MLP	30	0.8440971508067067	0.95	0.9515851711217844
LDA + MLP	30	0.0	1.0	1.0
(VGG + PCA) + LR	55	0.7905694150420949	0.9125	0.8916600529100529

## IIIT-CFW Dataset

For IIIT-CFW Dataset, the combination of ResNet with Logistic Regression worked best, giving an accuracy of 99.26% on the test set.

Classifier Method vs Accuracy Table for IIIT-CFW Dataset

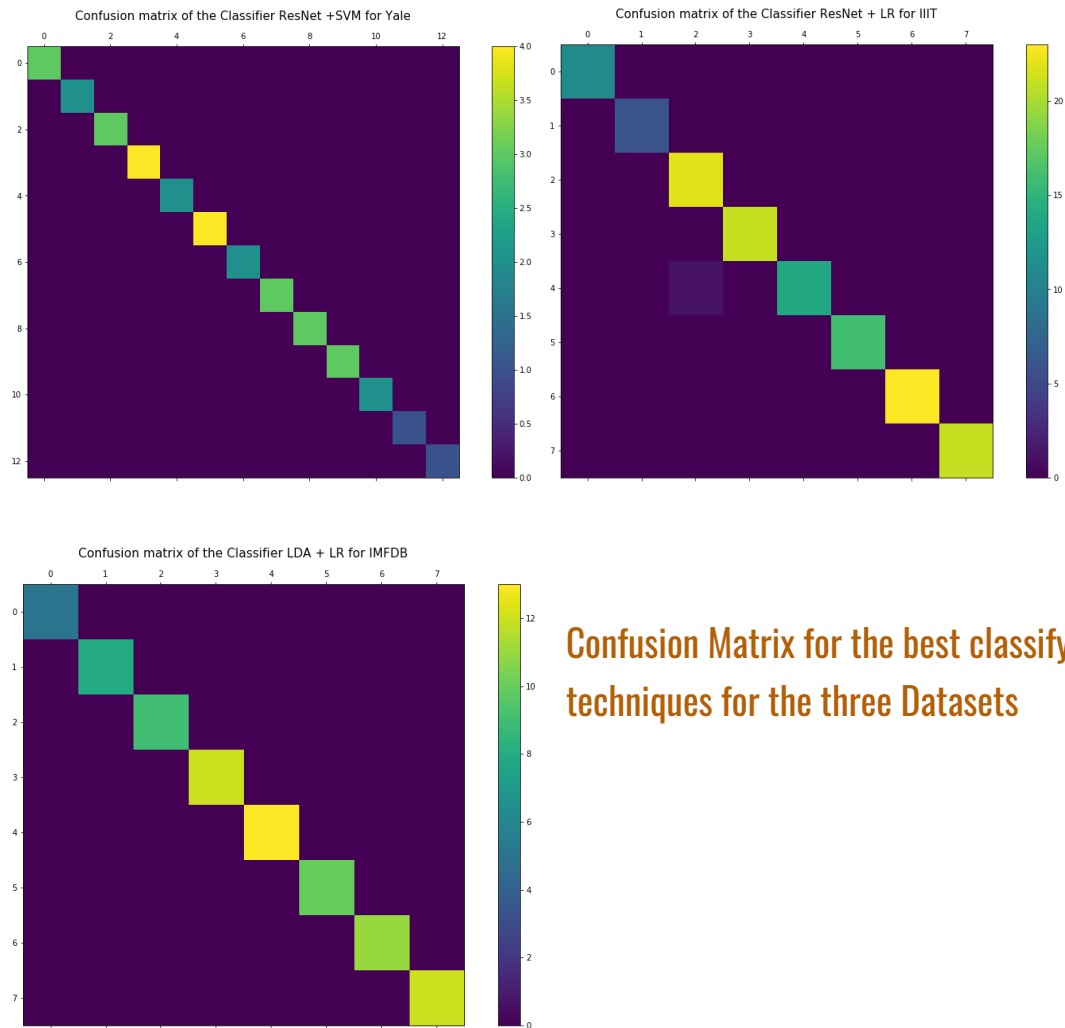
Method	Reduced Space	Classification Error	Accuracy	F1 Score
LDA + SVM	30	0.5163977794943222	0.9703703703703703	0.9627315513943421
ResNet + LR	2048	0.17213259316477408	0.9925925925925926	0.992911877394636
ResNet + SVM	2048	0.19245008972987526	0.9851851851851852	0.9853866740494648
(ResNet + KPCA) + LR	30	0.19245008972987526	0.9851851851851852	0.9853866740494648
(ResNet + KPCA) + MLP	30	0.5018484351393873	0.9629629629629629	0.9674560429787868
ResNet + SVM	2048	0.19245008972987526	0.9851851851851852	0.9853866740494648

## Yale Dataset

For Yale Dataset, the combination of Logistic Regression with all the three feature extraction + Resnet with SVM worked best, giving an accuracy of 100% on the test set.

Classifier Method vs Accuracy Table for Yale Dataset

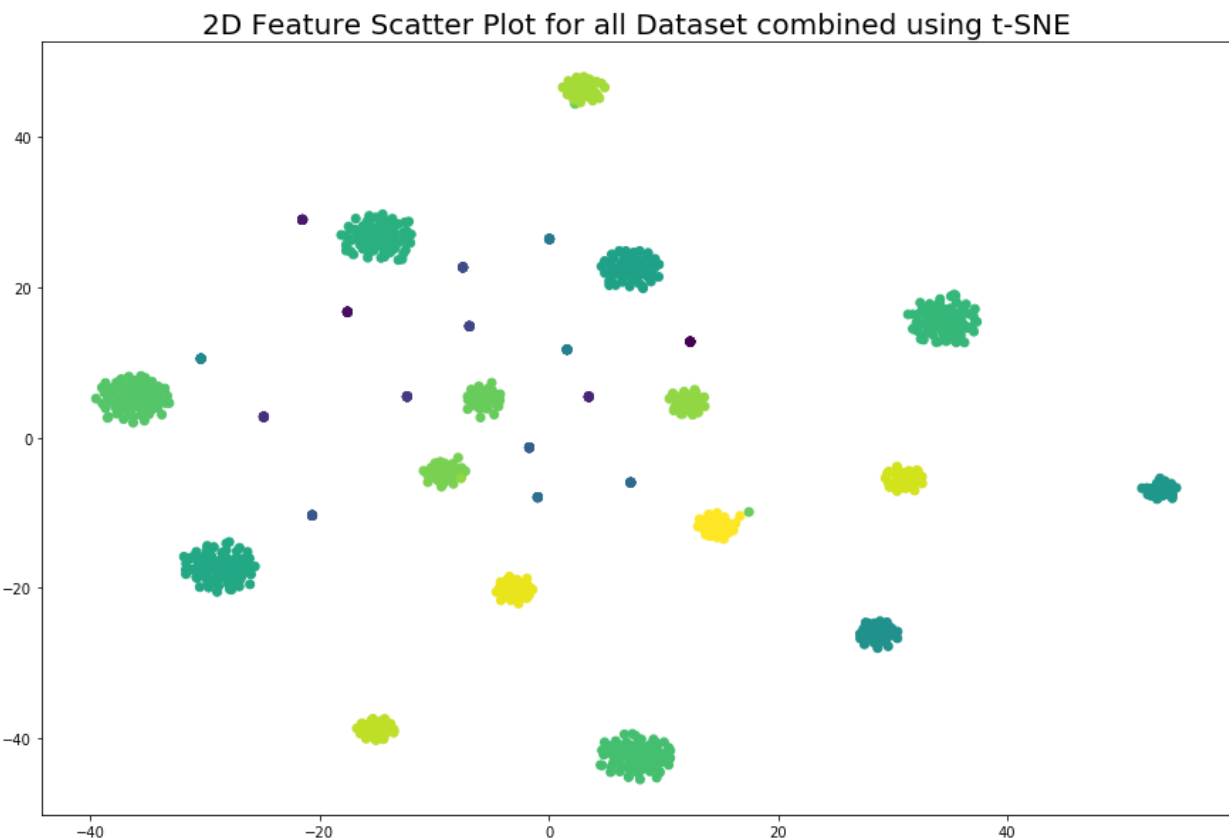
Method	Reduced Space	Classification Error	Accuracy	F1 Score
PCA + LR	55	2.3741027013091993	0.9090909090909091	0.85
LDA + LR	30	0.0	1.0	1.0
ResNet + SVM	2048	0.0	1.0	1.0
ResNet + LR	2048	0.0	1.0	1.0
(ResNet + KPCA) + LR	30	0.0	1.0	1.0
LDA + SVM	30	1.9462473604038075	0.9393939393939394	0.9040816326530612



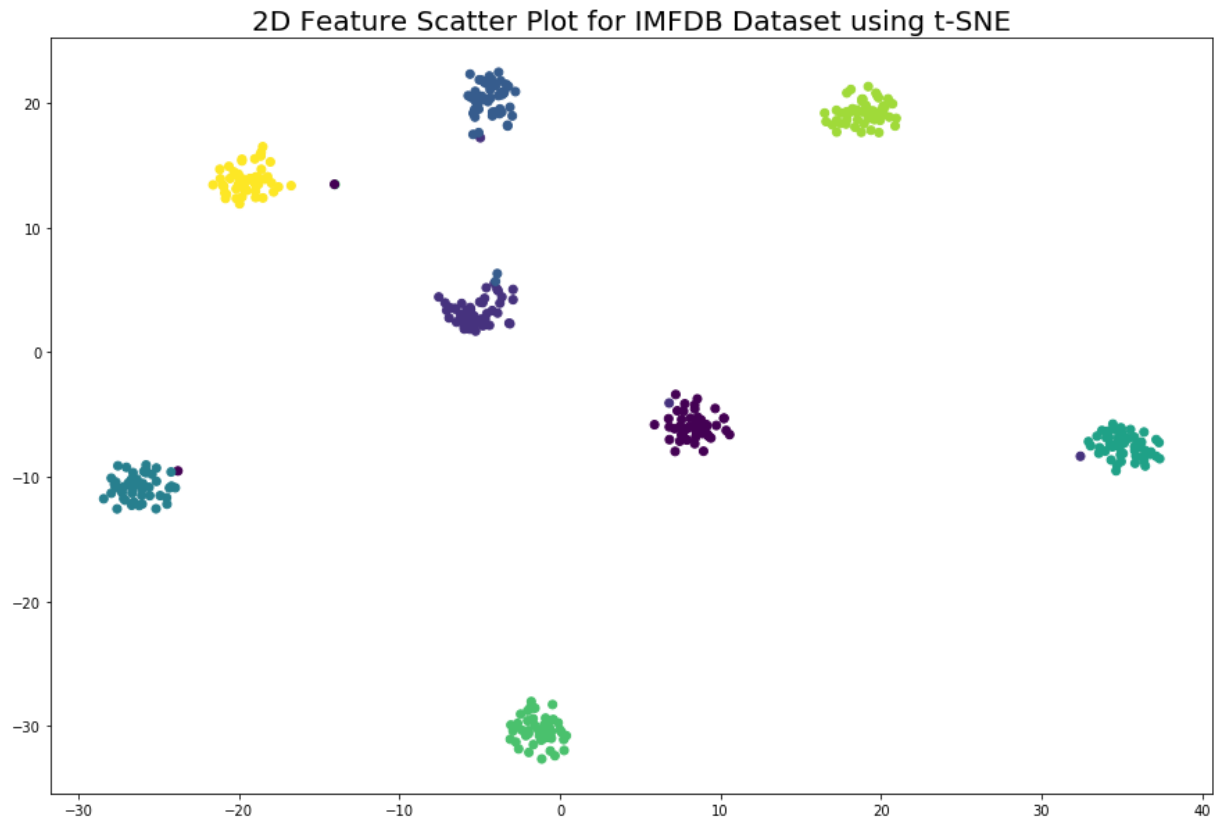
## 6. Does it make sense? Do you see similar people coming together? Can you do visualization dataset wise and combined?

t-Distributed Stochastic Neighbor Embedding (t-SNE) is an unsupervised, non-linear technique primarily used for data exploration and visualizing high-dimensional data.

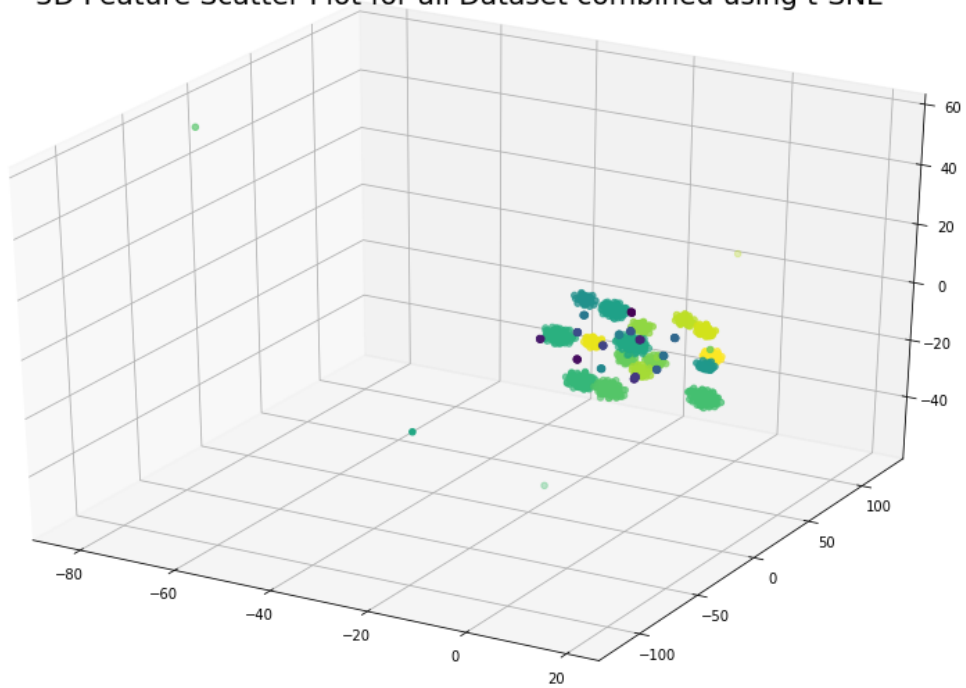
In the 2D/3D t-SNE scatter plot for all the datasets combined, it can be seen that samples with the same label are closer to each other than others with no exceptions.



In the 2D t-SNE scatter plot for only the IMFDB dataset also, it can be seen that samples with the same label are closer to each other forming bunches of similar samples.



### 3D Feature Scatter Plot for all Dataset combined using t-SNE





## 7. Classification with KNN

### a. How do we formulate the problem using KNN ?

We apply the above mentioned feature extraction methods to the dataset, then train the KNN classifier with the train set extracted from the complete dataset.

Given a sample and class id, we first apply the same feature extraction to the sample and then with the KNN, find the K-nearest neighbours to the sample and assign the label through majority voting. Then, we check if the,

given\_class\_id == assigned\_label, if so return True else False

### b. How do we analyze the performance ? Suggest the metrics (like accuracy) that is appropriate for this task.

Given a sample and KNN-classifier, we can use the accuracy metric as :-

$$\text{Accuracy} = \frac{\text{number of correct neighbours}}{\text{total neighbours}}$$

Or if we have multiple samples in the test set, we can use fraction of correctly classified test samples, MSE Loss or F-1 Score.

## KNN Classification with different Feature Extraction Method

Classifier Method vs Accuracy Table for IIIT-CFW Dataset

Method	Reduced Space	Classification Error	Accuracy	Precision Score
PCA + KNN	55	2.211083193570267	0.4666666666666667	0.4694764056606162
LDA + KNN	30	0.6776866969767514	0.9703703703703703	0.9739649936061381
ResNet + KNN	2048	0.3849001794597505	0.9555555555555556	0.9521554834054834
(ResNet + KPCA) + KNN	30	0.37515428924742517	0.9629629629629629	0.9571373674634545
(VGG + PCA) + KNN	55	2.023930902055773	0.6148148148148148	0.5663150042625746

Classifier Method vs Accuracy Table for IMFDB Dataset

Method	Reduced Space	Classification Error	Accuracy	Precision Score
PCA + KNN	55	2.2276669409945464	0.5625	0.559642094017094
LDA + KNN	30	0.0	1.0	1.0
ResNet + KNN	2048	0.6123724356957945	0.9625	0.9707792207792207
(ResNet + KPCA) + KNN	30	0.6123724356957945	0.9625	0.9707792207792207
(VGG + PCA) + KNN	55	1.0	0.875	0.845149642024642

Classifier Method vs Accuracy Table for Yale Dataset

Method	Reduced Space	Classification Error	Accuracy	Precision Score
PCA + KNN	55	2.8284271247461903	0.6666666666666666	0.6511111111111111
LDA + KNN	30	0.0	1.0	1.0
ResNet + KNN	2048	0.0	1.0	1.0
(ResNet + KPCA) + KNN	30	0.0	1.0	1.0
(VGG + PCA) + KNN	55	4.221158824088691	0.696969696969697	0.6666666666666665

## KNN Classification with LDA and different Nearest Neighbours

KNN with different Neighbour Number vs Accuracy Table for IIIT-CFW Dataset

Method	Neighbours	Reduced Space	Classification Error	Accuracy	Precision Score
LDA + KNN	1	30	0.5577733510227171	0.9703703703703703	0.9677891042780749
LDA + KNN	3	30	0.4472135954999579	0.9703703703703703	0.9687360739750446
LDA + KNN	5	30	0.7601169500660919	0.9629629629629629	0.967071611253197
LDA + KNN	9	30	0.7601169500660919	0.9629629629629629	0.967071611253197
LDA + KNN	12	30	0.7601169500660919	0.9629629629629629	0.967071611253197

KNN with different Neighbour Number vs Accuracy Table for IMFDB Dataset

Method	Neighbours	Reduced Space	Classification Error	Accuracy	Precision Score
LDA + KNN	1	30	0.11180339887498948	0.9875	0.9875
LDA + KNN	3	30	0.0	1.0	1.0
LDA + KNN	5	30	0.0	1.0	1.0
LDA + KNN	9	30	0.0	1.0	1.0
LDA + KNN	12	30	0.0	1.0	1.0

KNN with different Neighbour Number vs Accuracy Table for Yale Dataset

Method	Neighbours	Reduced Space	Classification Error	Accuracy	Precision Score
LDA + KNN	1	30	0.0	1.0	1.0
LDA + KNN	3	30	0.0	1.0	1.0
LDA + KNN	5	30	0.0	1.0	1.0
LDA + KNN	9	30	0.0	1.0	1.0
LDA + KNN	12	30	0.0	1.0	1.0

## 8. Gender Classification on (IMFDB + IIIT Dataset)

- a. We have a combination of IIIT-CFW and IMFDB datasets and we try to classify the face as male/female. Basically, the problems is a binary classification problem of where Male is assigned 0 and Female is assigned the label 1.

### Real Life Applications

- b. Gender recognition from face images is an important application in the fields of security, retail advertising and marketing.
- c. In places with crowd, if we can develop a model which can identify gender with satisfactory accuracy, then we can keep an easy track of moving crowd.
- d. It could prove useful at railway stations, where opposite gender might try to avail benefits based on gender.

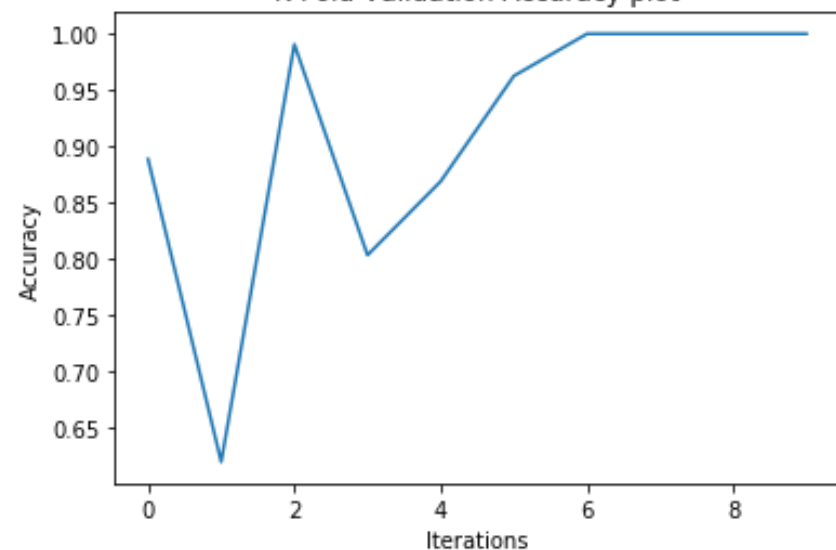
### Model Pipeline

- e. Load the IMFDB and IIIT-CFW Dataset, and concatenate the samples from both datasets to create single dataset.
- f. Change the labels for Male class as 0 and for Female class as 1.
- g. Use LDA Feature Extraction as it works good(seen in previous classifications)
- h. Apply all the classifier models like SVM, Logistic Regression, MLP and KNN for training (Best results by SVM and KNN ~98% accuracy)

- i. Use different accuracy metrics like **Accuracy, MSE, Precision, F1-score** to measure the correctness of the model

## K-Fold Validation

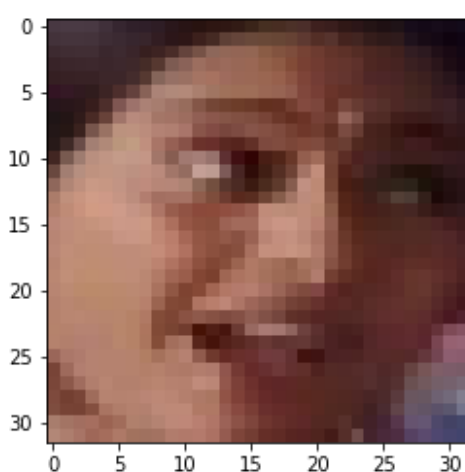
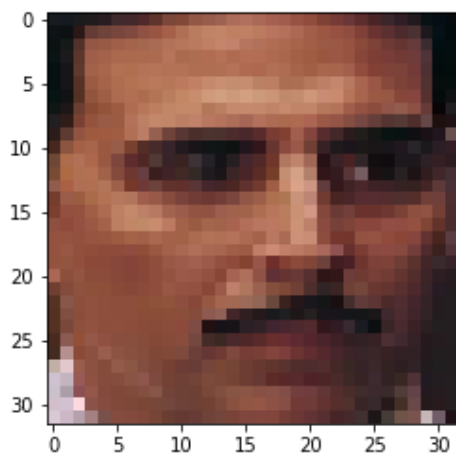
K-Fold Validation Accuracy plot



Accuracy Table for different classifiers for Gender Prediction

Classifiers	Accuracy
SVM	0.986046511627907
Logistic Regression	0.9813953488372092
MLP	0.9813953488372092
KNN	0.986046511627907

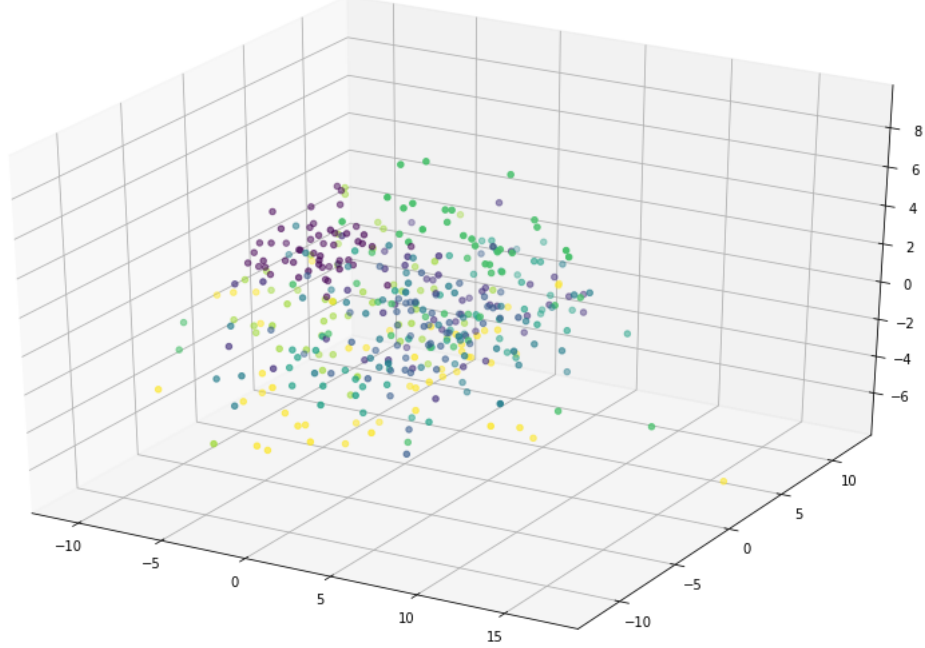
## Few Misclassified Images



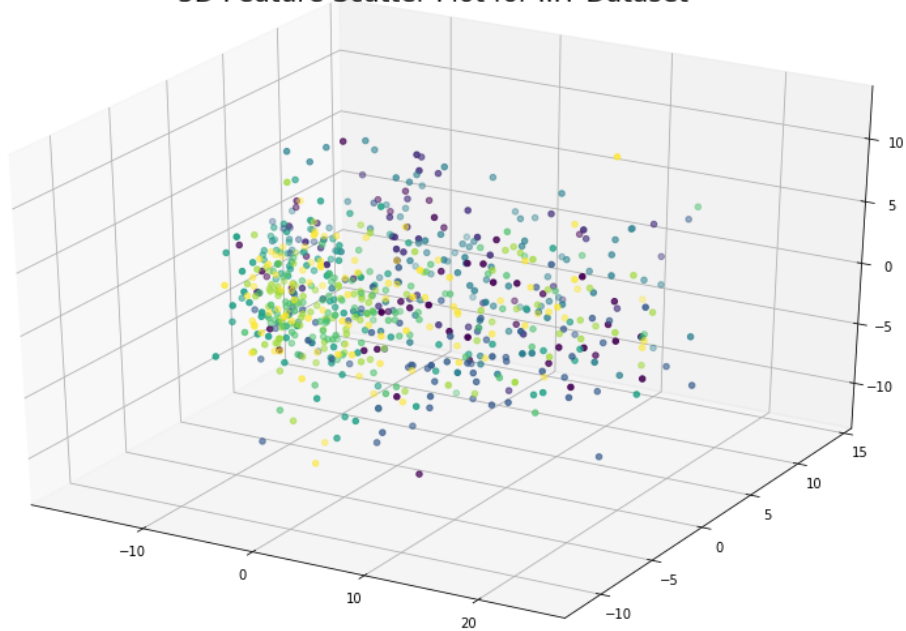
## Other Plots :-

### PCA

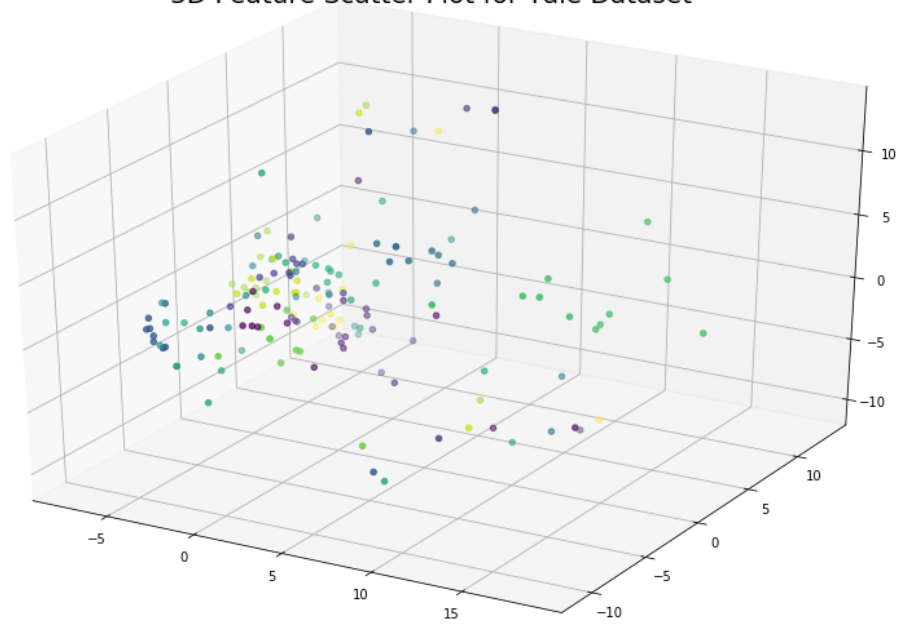
3D Feature Scatter Plot for IMFDB Dataset



3D Feature Scatter Plot for IIIT Dataset



3D Feature Scatter Plot for Yale Dataset



t-SNE

3D Feature Scatter Plot for IMFDB Dataset using t-SNE

