

Implementing Various Classifiers to Achieve Facial Recognition

Omar H. Abdelkader¹

¹University of Maryland, College Park
Electrical and Computer Engineering – A.V. Williams Building
8223 Paint Branch Drive – College Park, MD – 20740

oabdelka@umd.edu

Abstract. *The purpose of this technical report is to analyze the performance of various classifier implementations in the context of facial recognition. Specifically, the Bayesian and K-nearest neighbors classifiers are considered. In addition, Principal Component Analysis (PCA) and Fisher's Multiple Discriminant Analysis (MDA) are evaluated as dimensionality reduction techniques. The aforementioned methods have been tested on three different data sets, including images with various pose, expression, and illumination variations.*

Keywords: *Bayesian decision theory, K-nearest neighbors, PCA, MDA.*

1. Introduction

Thus far, we have covered two classification methodologies. In the *bottom-up* approach, the training data is first fit to some known distribution using a procedure such as maximum likelihood estimation. The discriminant functions are computed from the distribution, to which the testing data can then be classified. The Bayesian classifier, otherwise known as the minimum error rate classifier, belongs to this category. Conversely, the *top-down* approach begins with choosing a classifier geometry, and aims to fit the data to that geometry afterwards. The top-down approach lends itself nicely to data that does not appear to fit a known distribution, but still has an intuitive and clear separation between the classes. The K-Nearest Neighbor Rule is a non-parametric classifier that is often used in these cases.

In addition to classification techniques, we have also examined the importance of dimensionality reduction. In some cases, many features can be overlooked while still maintaining the accuracy of a particular classifier. Dimensionality reduction techniques are vital for both reducing computation time and reducing the number of samples required to create a suitable classifier. The first technique covered in this course was Principal Component Analysis. PCA works to represent the given data in a lower dimension space such that the representation captures as much energy of the data as possible. In some cases, though, projecting onto a lower dimension space produced by PCA destroys the class separation. Fisher's Multiple Discriminant Analysis aims to fix this flaw by maximizing between class scatter in conjunction with minimizing within class scatter.

2. Theoretical Framework

In this project, I examined three `.mat` data sets. The *face* data set is composed of cropped images of 200 subjects, three images each. The first image is a neutral face, the second

has a facial expression, and the third has illumination variations. The *pose* data set is composed of cropped images of 68 subjects under 13 different poses, and the *illumination* data set is composed of cropped images of 68 subjects under different 21 illuminations.

Table 1. Data set Breakdown

	<i>Num. of Classes</i>	<i>Samples per Class</i>	<i>Image Size</i>
Face	200	3	24 x 21
Pose	68	13	48 x 40
Illumination	68	21	48 x 40

The number of subjects in each data set can be thought of as the number of classes. In addition, since each pixel is considered to be a unique feature, the image size directly corresponds to the number of dimensions for the given model. Two photos from the *pose* data set are presented below. The images are of the same subject, but under different poses.



Figure 1. Sample 1, Pose 1



Figure 2. Sample 1, Pose 2

3. Methodology

Prior to conducting any analysis or manipulation of the data, I first reshaped the image matrices into feature vectors. Fortunately, this is trivial to do in MATLAB, and the *illumination* data set is already provided in vector form. Once reshaped, each data set was divided into two disjoint sets, one for training, and another for testing. For the purposes of this paper, I divided the *face* data set into $\frac{2}{3}$ training to $\frac{1}{3}$ testing, and the *pose* and *illumination* data sets to $\frac{3}{4}$ training to $\frac{1}{4}$ testing.

3.1. Bayes' Classifier

In order to correctly use the Bayes' classifier, the data must first be fit to a known distribution. Fortunately, in this project, we are told that the underlying distribution of the images is Gaussian. Estimating the Gaussian parameters $\hat{\mu}$ and $\hat{\Sigma}$ is trivial, as we have done the maximum likelihood estimation derivation for the multivariate Gaussian in class.

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n \vec{x}_k \quad (1)$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^n (\vec{x}_k - \hat{\mu})(\vec{x}_k - \hat{\mu})^T \quad (2)$$

With $\hat{\mu}$ and $\hat{\Sigma}$ from equations 1 and 2, computing the posterior probability should be simple. Since the prior and evidence probabilities are equal across all classes, it is safe to compute and classify using class conditional probability values.

$$P(\omega_i|\vec{x}) = \frac{1}{(2\pi)^{\frac{d}{2}}|\hat{\Sigma}_i|^{\frac{1}{2}}} e^{-\frac{1}{2}(\vec{x}-\hat{\mu}_i)^T\hat{\Sigma}_i^{-1}(\vec{x}-\hat{\mu}_i)} \quad (3)$$

In general, any positive monotonic transformation of the discriminant function will classify exactly the same as the discriminant itself. So, to simplify even further, I can take the natural log of the class conditional probability and rearrange the terms. The boxed term in equation 4 below can be dropped because it is constant for all i .

$$\begin{aligned} g_i(\vec{x}) &= \ln P(\omega_i|\vec{x}) \\ &= -\frac{1}{2}(\vec{x} - \hat{\mu}_i)^T \hat{\Sigma}_i^{-1} (\vec{x} - \hat{\mu}_i) - \boxed{\frac{d}{2} \ln 2\pi} - \frac{1}{2} \ln |\hat{\Sigma}_i| \\ &= -\frac{1}{2}(\vec{x}^T \hat{\Sigma}_i^{-1} \vec{x} - 2\hat{\mu}_i^T \hat{\Sigma}_i^{-1} \vec{x} + \hat{\mu}_i^T \hat{\Sigma}_i^{-1} \hat{\mu}_i) - \frac{1}{2} \ln |\hat{\Sigma}_i| \\ &= \vec{x}^T (-\frac{1}{2} \hat{\Sigma}_i^{-1}) \vec{x} + \hat{\mu}_i^T \hat{\Sigma}_i^{-1} \vec{x} + (-\frac{1}{2} \hat{\mu}_i^T \hat{\Sigma}_i^{-1} \hat{\mu}_i - \frac{1}{2} \ln |\hat{\Sigma}_i|) \\ &= \vec{x}^T W_i \vec{x} + w_i^T \vec{x} + w_{0i} \end{aligned} \quad (4)$$

$$(5)$$

In this case, however, computing $g_i(\vec{x})$ is not as trivial plugging in and evaluating equation 5. For any $\vec{x} \in \mathbb{R}^d$, where d is greater than the number of samples for a given class, n , the estimate of the covariance matrix $\hat{\Sigma}$ will always be singular. This means $\hat{\Sigma}$ is non-invertible and has a determinant of zero.

Examining the procedure for computing $\hat{\Sigma}$ reveals why this happens. The multiplication of a vector with its transpose will always be a matrix of rank one because the columns of the resulting matrix are scaled versions of the original vector. For some vector, \vec{x}

$$\vec{x} \cdot \vec{x}^T = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_d \end{bmatrix} = \begin{bmatrix} x_1x_1 & x_2x_1 & \dots & x_dx_1 \\ x_1x_2 & x_2x_2 & \dots & x_dx_2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1x_d & x_2x_d & \dots & x_dx_d \end{bmatrix}$$

Adding n matrices of this form will increase the rank one by one up to n , however if $n < d$, there will not be enough pivot columns in the end $d \times d$ matrix. Taking the inverse of $\hat{\Sigma}$ is what is known as an *ill-posed* inference problem. A common fix for regularizing problems of this kind is to use shrinkage estimation, whereby a naive or raw estimate is improved by combining it with other information. In this case, $\hat{\Sigma}$ becomes positive-definite by adding a small value to all the diagonal elements, namely the identity matrix.

3.2. K-Nearest Neighbors

Unlike the Bayesian classifier, the K-Nearest Neighbors classifier is non-parametric. The algorithm for classifying according to the K-nearest neighbors rule is as follows

Data: Gather n labeled samples $\{x_1, \dots, x_n\}$ and choose a tiebreaker mode
for a new sample \vec{x} **do**
 Find k -nearest neighbors, $\vec{x}'_1, \dots, \vec{x}'_k \in \{\vec{x}_1, \dots, \vec{x}_n\}$, closest to \vec{x} ;
 Assign \vec{x} label of majority of it's k -nearest neighbors;
 if multiple k -nearest neighbors **then**
 Use the tiebreaking mode provided;
 end
end

Algorithm 1: K-Nearest Neighbors Algorithm

In the event of a tie (i.e. two or more classes have the same number of closest neighbors to the sample), a tiebreaking mode of the user's choice is used to classify. The following tiebreaking methods are investigated in this report.

Table 2. K-Nearest Neighbors Tiebreaking Methods

<i>Tiebreaker Mode</i>	<i>Description</i>
Discard	Samples with a tie are discarded
Retry	Repeat with $K - 1$ neighbors until the tie is broken
Closest	Of the tied classes, select the one with the minimum total distance
Random	Of the tied classes, select one randomly

3.3. Principal Component Analysis

Principal Component Analysis is a technique for reducing a given data set to a lower dimension. PCA works by projecting the original data onto a subset of the feature vectors, known as the principal components. In order to compute which directions to project onto, one computes the scatter matrix of the centered data and solves for its eigenvalues and eigenvectors. By projecting onto the eigenvectors corresponding to the largest eigenvalues, we are guaranteed to maintain as much of the original data's variance as possible. The extent to which how much data is eliminated is determined by how much energy the operator is willing to sacrifice. This value is expressed by the parameter, α .

For example, let d be the number of features for a particular data set. Let m be the number of features we wish to reduce our data to, where $m < d$. Then, the energy preserved by projecting onto the eigenvectors corresponding to the m largest eigenvalues can be expressed as R_m

$$\begin{aligned}
 R_m &= \frac{\sum_{i=1}^m \lambda_i}{\sum_{i=1}^d \lambda_i} \\
 &= \frac{\lambda_1 + \lambda_2 + \dots + \lambda_m}{\lambda_1 + \lambda_2 + \dots + \lambda_d} = 1 - \alpha
 \end{aligned} \tag{6}$$

It is important to note that the correct way to use PCA is to run it on the training set, save resulting the principal components, and then use them to transform the testing set. This

is the only way to guarantee that both sets end up in the same space without using any knowledge about the test set during training.

3.4. Fisher’s Multiple Discriminant Analysis

A major shortcoming of PCA is that it does not preserve between-class scatter very well. In other words, samples that appear to belong to distinct classes in a higher dimension space appear to belong to the same class when projected onto a lower dimension space. Fisher’s Linear Discriminant Analysis projects a 2-class data set of d dimensions onto a line for which the projected samples are well separated.

Generalizing this procedure to c classes is known as Multiple Discriminant Analysis (MDA). Naturally, this involves at most $c - 1$ discriminant functions; thus, the projection is from a d -dimensional space to an m -dimensional space, where $m \leq c - 1 < d$. We seek a transformation matrix \mathbf{W} that maximizes the ratio of the between-class scatter to the within-class scatter. Using this measure, we define a criterion function, $J(\mathbf{W})$

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|} \quad (7)$$

The columns of an optimal \mathbf{W} that maximize equation 7 are the generalized eigenvectors corresponding to the m -largest eigenvalues in

$$\mathbf{S}_B \mathbf{w}_i = \lambda_i \mathbf{S}_W \mathbf{w}_i \quad (8)$$

Just as in PCA, MDA should be run solely on the training data. The transformation matrix that is generated can then be used to reduce the testing data to m dimensions.

4. Results

I divided each of the three data sets into the following training and testing sets according to the ratios defined in Section 3.

Table 3. Training-Testing Set Split

	<i>Training Samples</i>	<i>Testing Samples</i>
Face	{1, 2}	{3}
Pose	{1-10}	{11-13}
Illumination	{1-16}	{17-21}

4.1. Bayes’ Classifier

As expected, out of the three data sets considered, the Bayesian classifier performed poorest on the *face* data set. Due to the fact that there were only two training samples available for every class, I did not have very high expectations for any of the classifiers to perform well on this data set. Despite this, it is quite astonishing that the Bayesian classifier managed to correctly label 64% of the testing data. On average, randomly assigning a label

of the 200 possible classes would have yielded 0.5% accuracy. The fact that the Bayesian classifier performed 128 times better than randomly guessing on a data set with only two training samples per class really speaks to the power of its classification ability.

Table 4. Bayesian Classifier Performance

	<i>Correct Classifications</i>	<i>Misclassifications</i>	<i>Accuracy</i>
Face	128	72	64%
Pose	139	65	68.1%
Illumination	340	0	100%

In spite of having 5 times as many training samples per class, the Bayesian classifier performed nearly just as poorly on the *pose* data set as on the *face* data set. On the other hand, the Bayesian classifier correctly classified all of the testing samples in the *illumination* data set.

According to these results, the Bayesian classifier is certainly illumination invariant, but the same cannot be said for pose variations. Perhaps a data set with more samples would yield results more indicative of the classifier's performance. The covariance matrices in all three data sets had to be regularized in order to be used in any computation, meaning that the data was not well-formed for this kind of problem.

4.2. K-Nearest Neighbors

For all three data sets, the K-Nearest Neighbors classifier performed equal-to, or worse than the Bayesian classifier. For $K = 1$, the *face* and *pose* data sets were 4.5% and 6.3% less accurate compared to the Bayesian classifier, respectively. The nearest neighbor classifier, however, was just as accurate for the *illumination* data set, correctly classifying every sample in the test set.

Due to the limited number of samples, both per class and overall, I only varied K from one to five. Intuitively, it would make sense that corroborating the label assignment against more neighbors from the training set would improve the labeling accuracy of the model. Surprisingly, however, increasing K seemed to hurt the overall precision of the model.

Despite the Bayesian classifier being more accurate across all the data sets, it is worth mentioning the price paid in execution time to achieve the modest increase in accuracy. Take the *illumination* testing set, for example. It took 293.87 seconds to compute the Bayesian classifier predictions, whereas it took only 4.94 seconds to compute the K-Nearest Neighbor predictions ($K = 1$). For high dimension data, taking the inverse and determinant of the covariance matrix can be quite expensive. As a result, sacrificing a few percentage points in accuracy could be worth using a simpler classifier such as K-Nearest Neighbors.

4.2.1. Tiebreaker Strategies

The `discard` tiebreaker strategy is not a viable classification technique in practice. This is because samples in the testing set with numerous nearest neighbors will wind up unlabeled.

beled. As a result, when overlooking unclassified samples, this adaptation of the nearest neighbors has high accuracy. It performs rather poorly, however, if the unclassified samples are counted as misclassifications. In the tables below, I chose to report the number of discarded samples as opposed to the accuracy of the `discard` method.

The `retry` and `closest` tiebreaking implementations produced nearly identical results of one another. Samples that tie in the $K = 2$ variant are always classified as they would have been had $K = 1$. As a result, the $K = 2$ accuracy is equivalent to the $K = 1$ accuracy for both tiebreaking modes. In addition, these modes saw the least precipitous drop in accuracy as K is increased. This outcome is likely due to the fact that these two strategies are the only ones that use information about the initial run’s nearest neighbors when attempting to reclassify.

Table 5. Face Data set Results

Discard				Retry			
	<i>Correct</i>	<i>Incorrect</i>	<i>Discarded</i>		<i>Correct</i>	<i>Incorrect</i>	<i>Accuracy</i>
K = 1	119	81	0	K = 1	119	81	59.5%
K = 2	53	3	144	K = 2	119	81	59.5%
K = 3	67	8	125	K = 3	123	77	61.5%
K = 4	68	19	113	K = 4	120	80	60%
K = 5	69	28	103	K = 5	119	81	59.5%
Closest				Random			
	<i>Correct</i>	<i>Incorrect</i>	<i>Accuracy</i>		<i>Correct</i>	<i>Incorrect</i>	<i>Accuracy</i>
K = 1	119	81	59.5%	K = 1	119	81	59.5%
K = 2	119	81	59.5%	K = 2	98	102	49%
K = 3	123	77	61.5%	K = 3	94	106	47%
K = 4	120	80	60%	K = 4	91	109	45.5%
K = 5	119	81	59.5%	K = 5	91	109	45.5%

The final tiebreaking strategy, `random`, was initially implemented in anticipation of the other strategies requiring much more execution time to arrive at a result. As it so happens, however, benchmarking the strategies showed no discernible reduction in wait time when using `random` over `retry` or `closest`. Of course, it is important to note that this behavior can be partially attributed to the unusually small size of the data sets being used.

In practice, where the goal is to correctly classify each sample in the testing set as best as possible, a strategy like `retry` or `closest` would most likely be used to break a tie. Because `closest` outperformed `retry` in my results, albeit somewhat marginally, I chose to use `closest` for the remaining computations pertaining to K-Nearest-Neighbors that do not explicitly state a specific tiebreaker mode for the rest of the report.

Table 6. Pose Data set Results

Discard				Retry			
	<i>Correct</i>	<i>Incorrect</i>	<i>Discarded</i>		<i>Correct</i>	<i>Incorrect</i>	<i>Accuracy</i>
K = 1	126	78	0	K = 1	126	78	61.8%
K = 2	29	4	171	K = 2	126	78	61.8%
K = 3	46	22	136	K = 3	120	84	58.8%
K = 4	48	40	116	K = 4	114	90	55.9%
K = 5	56	50	98	K = 5	110	94	53.9%
Closest				Random			
	<i>Correct</i>	<i>Incorrect</i>	<i>Accuracy</i>		<i>Correct</i>	<i>Incorrect</i>	<i>Accuracy</i>
K = 1	126	78	61.8%	K = 1	126	78	61.8%
K = 2	126	78	61.8%	K = 2	87	117	42.7%
K = 3	120	84	58.8%	K = 3	74	130	36.3%
K = 4	115	89	56.3%	K = 4	71	133	34.8%
K = 5	111	93	54.4%	K = 5	72	132	35.3%

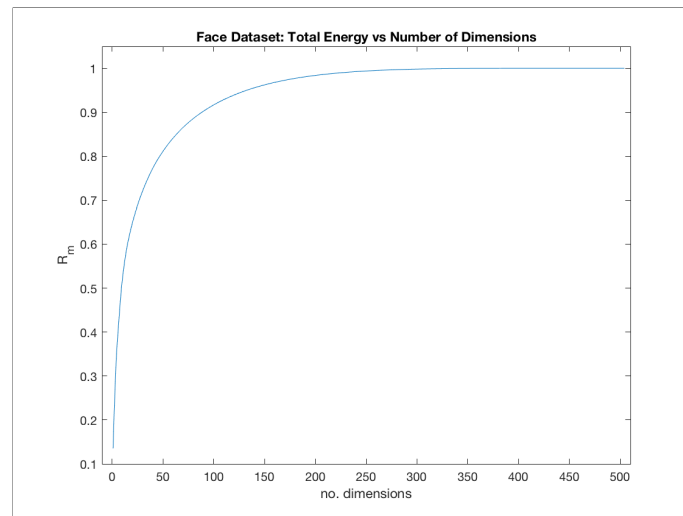
Table 7. Illumination Data set Results

Discard				Retry			
	<i>Correct</i>	<i>Incorrect</i>	<i>Discarded</i>		<i>Correct</i>	<i>Incorrect</i>	<i>Accuracy</i>
K = 1	340	0	0	K = 1	340	0	100%
K = 2	287	0	53	K = 2	340	0	100%
K = 3	290	7	43	K = 3	333	7	97.9%
K = 4	265	11	64	K = 4	327	13	96.2%
K = 5	250	17	73	K = 5	316	24	92.9%
Closest				Random			
	<i>Correct</i>	<i>Incorrect</i>	<i>Accuracy</i>		<i>Correct</i>	<i>Incorrect</i>	<i>Accuracy</i>
K = 1	340	0	100%	K = 1	340	0	100%
K = 2	340	0	100%	K = 2	316	24	92.9%
K = 3	333	7	97.9%	K = 3	304	36	89.4%
K = 4	329	11	96.8%	K = 4	283	57	83.2%
K = 5	322	18	94.7%	K = 5	280	60	82.4%

4.3. Principal Component Analysis

In the graph below, I plot R_m as a function of the number of retained features in the *face* data set. The images from this data set are originally of size 24 x 21 pixels, meaning that the feature vector, \vec{x} , is of size 504. Recall, PCA can be used to reduce the feature vector size to any arbitrary size m , so long as $m < d$, where $\vec{x} \in \mathbb{R}^d$.

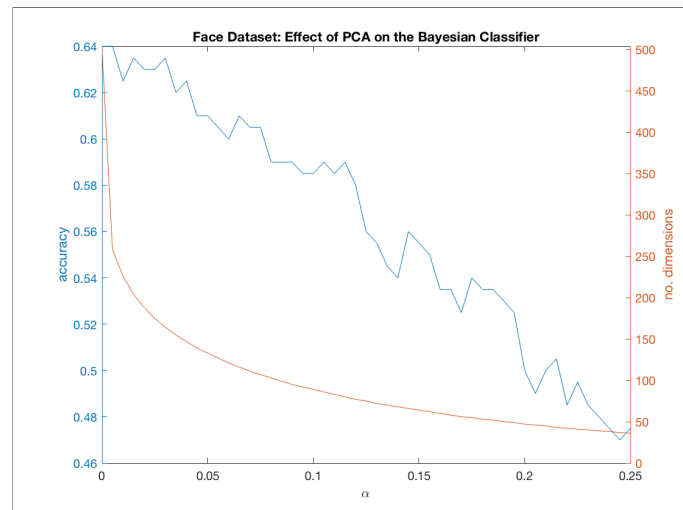
For example, reducing the *face* feature set in half (onto the 252 largest-variance principal components) would only cost 0.57% of the original energy. In fact, projecting onto 89 principal components would only lose 10% of the original energy.



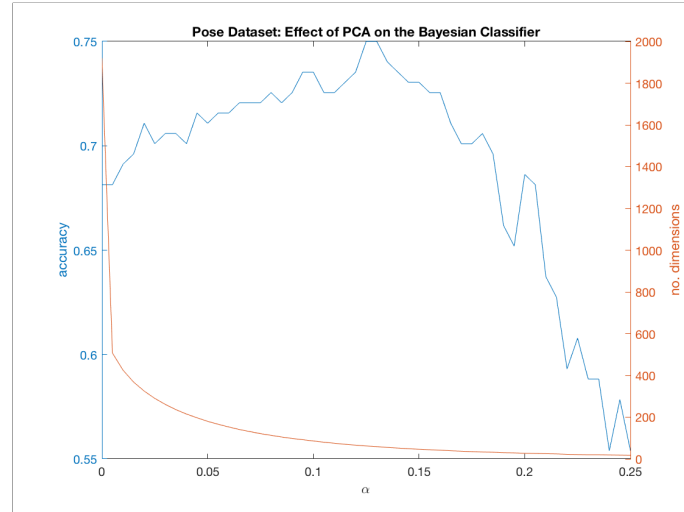
This same procedure can be used to reduce the feature set of the *pose* and *illumination* data sets.

4.3.1. Post-PCA Bayes'

The following graphs showcase the performance of the post-PCA Bayesian classifier as a function α .

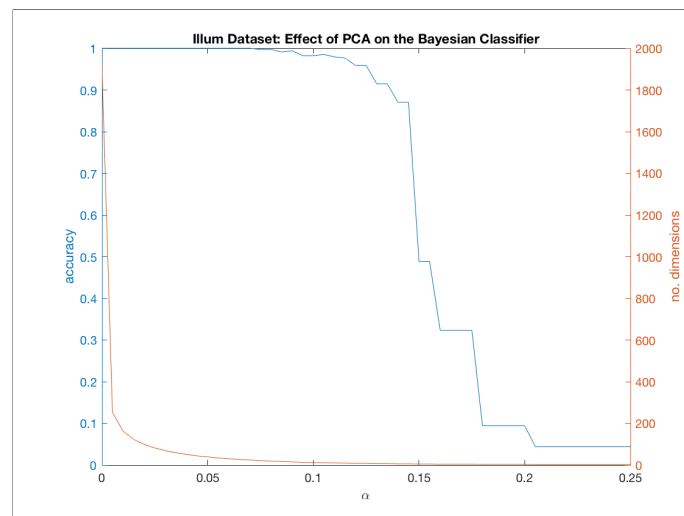


Projecting the *face* data set onto the 133 highest-variant dimensions ($\alpha = 0.05$) resulted in a loss of accuracy amounting to 3%. This drop-off in accuracy directly correlates to the low number of training samples used in this data set.



Conversely, PCA was much more effective on the *pose* and *illumination* data sets. The curves denoting the number of dimensions fell off much more quickly than on the *face* data set, presumably due to the number of training samples available. Using $\alpha = 0.05$ reduced the feature vector size of the *pose* data set from 1920 to 180, improving the accuracy of the model by 2.9%. For the illumination data set, $\alpha = 0.05$ reduced the feature vector size from 1920 to 39, at no cost to the accuracy of the model.

In other words, the Bayesian classifier still correctly classified all of the testing samples using while only 2% of the original features. Recall, the standard Bayesian classifier took 293.87 seconds to classify every test sample. After reducing the testing samples using PCA, this same feat only took 0.2049 seconds.



It would appear that many of the features in the *illumination* data set are extremely redundant. Fortunately, as one can see from the sample photos immediately on the following page, this checks out. Many of the pixels in these photos covary with the illumination alterations in the same manner (i.e. the pixels carry redundant information.)



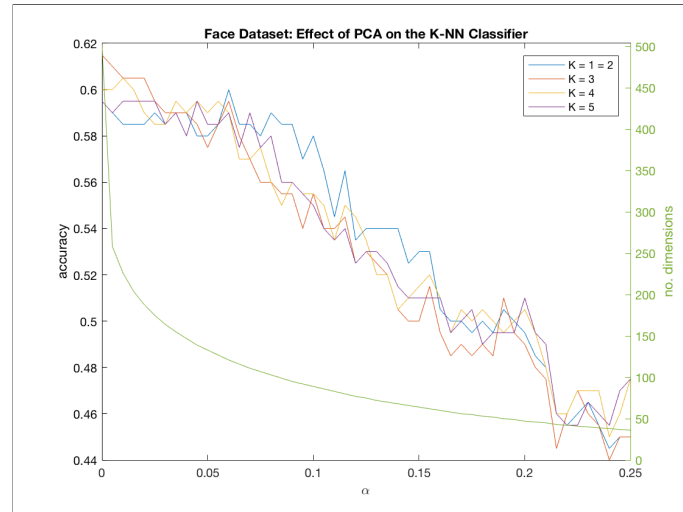
Figure 3. Sample 1, Illumination 1



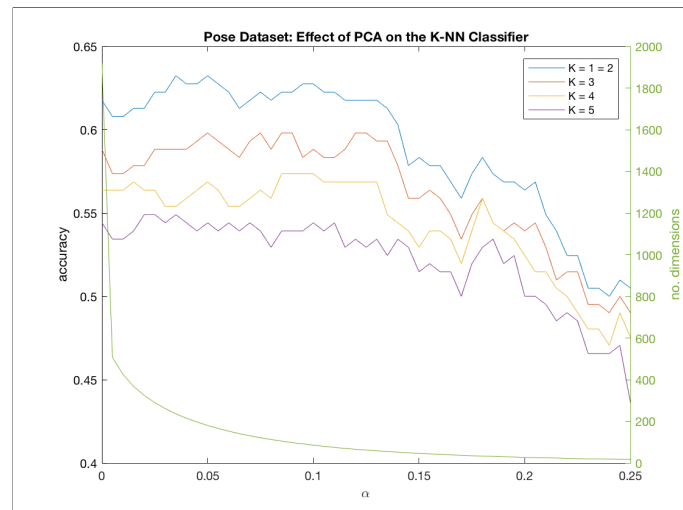
Figure 4. Sample 1, Illumination 5

4.3.2. Post-PCA K-Nearest Neighbors

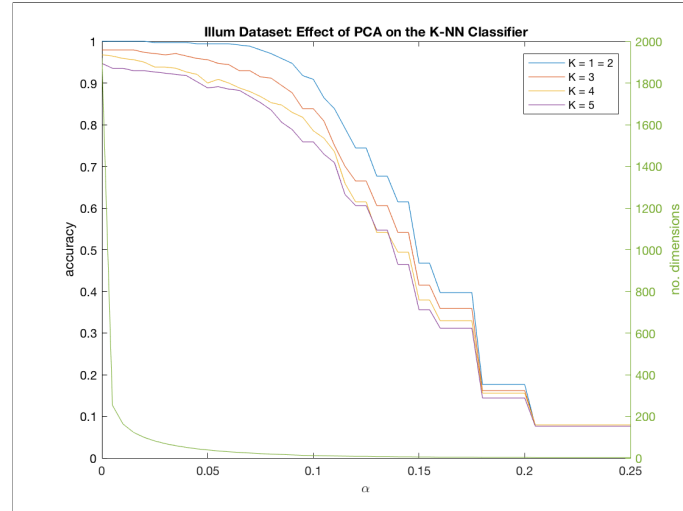
This section investigates the effect of varying the PCA parameter, α , on the K-Nearest Neighbors classification performance. As one can tell rather quickly, the shapes of the graphs for each particular data set bear a similar resemblance to their post-PCA Bayesian classifier counterparts.



Across all three data sets, the trend of higher K values performing worse also continued. As before, this is most likely because of the sparse number of samples in the training set.



Interestingly, the discrepancy in performance between the different K -variants seemed to converge as α increased. This makes sense due to the shrinkage of the class separation that occurs when utilizing PCA.

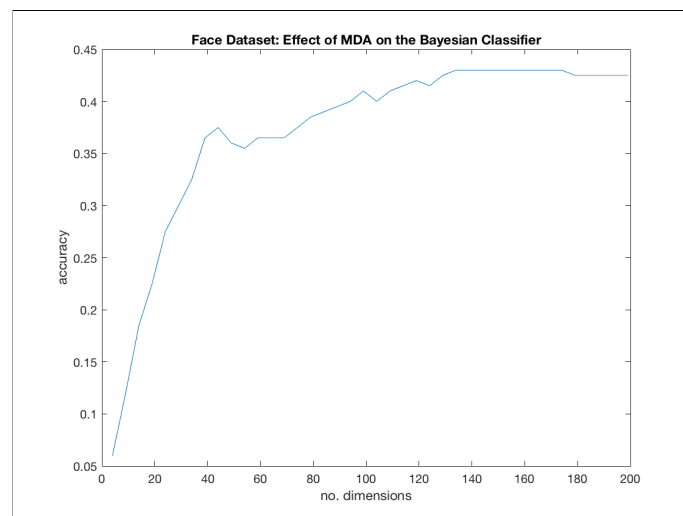


4.4. Fisher's Multiple Discriminant Analysis

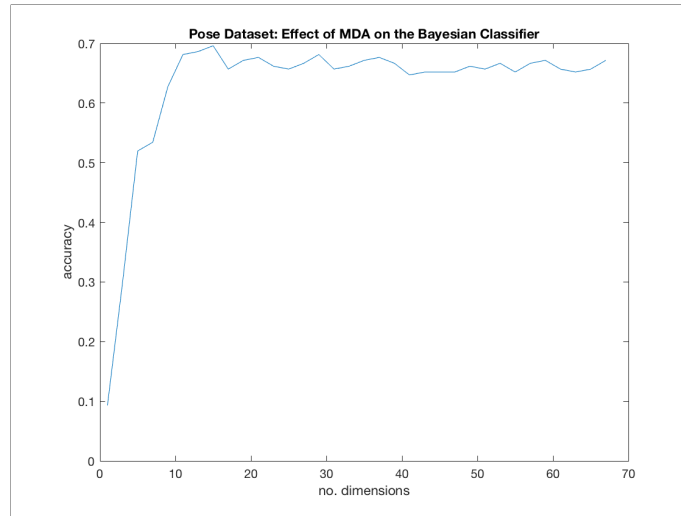
When used as a dimensionality reduction routine, Fisher's Multiple Discriminant Analysis projects any given data set onto at most $c - 1$ dimensions, where c is the number of classes. For the *face* data set, this requires projecting from 504 dimensions to at most 199. And for the others, it involves projecting from 1920 dimensions to at most 67.

4.4.1. Post-MDA Bayes'

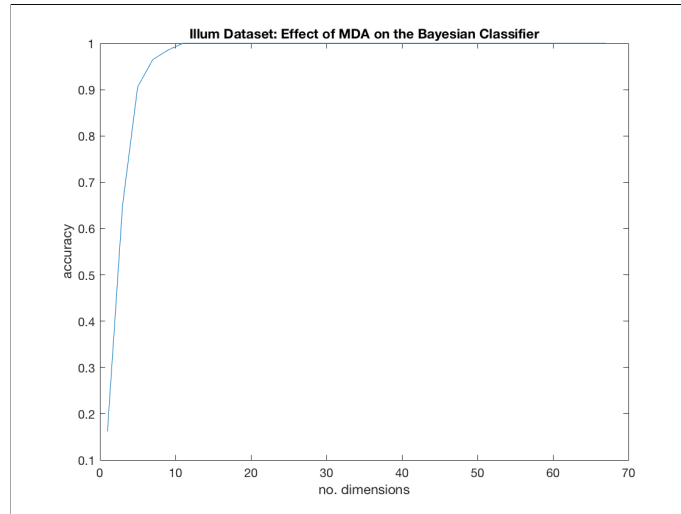
The following graphs showcase the performance of the post-MDA Bayesian classifier as a function of the number of dimensions the data is projected onto.



The performance of the post-MDA Bayesian classifier on the *face* data set was very poor. The accuracy was nearly 20% worse than prior to the MDA transformation. This behavior is very similar to that of the post-PCA transformation on the same data set.



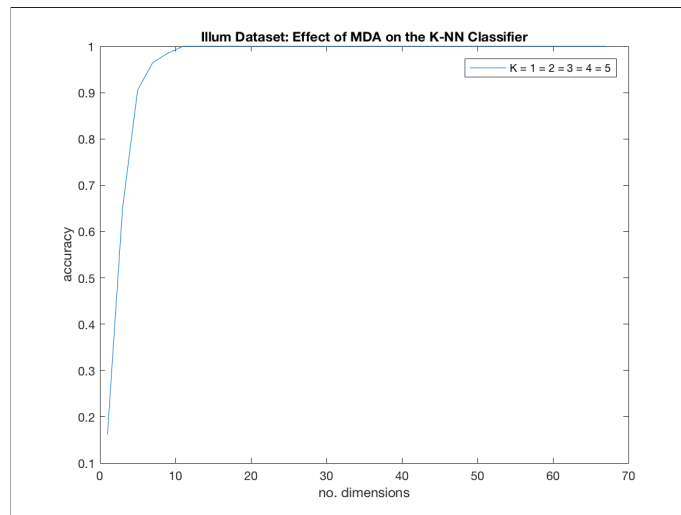
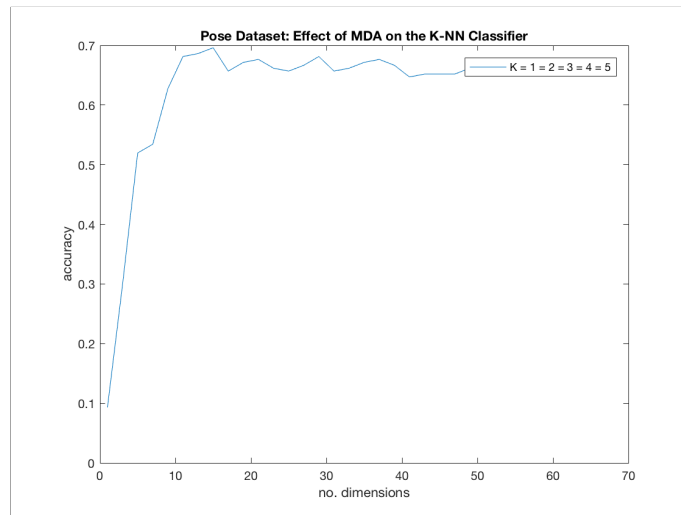
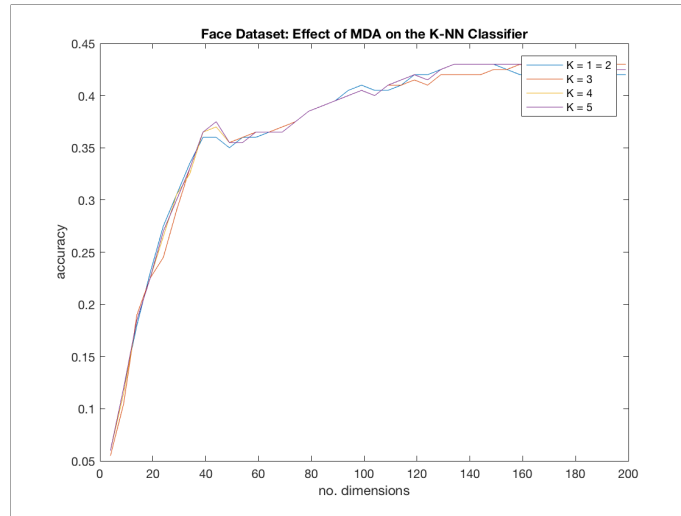
On the other hand, the post-MDA Bayesian classifier was excellent. Both the *pose* and *illumination* data sets, could be reduced to as low as 10 to 15 dimensions at no cost to the accuracy of the model. In fact, for some values of the dimension size, the Bayesian classification accuracy on the *pose* data set improved.



When comparing these plots to the PCA results, it is immediately apparent how important the between class scatter is to maintaining a low error rate. The Bayesian classifier performed better in a lower dimension space when reduced using MDA than in higher dimension spaces using PCA.

4.4.2. Post-MDA K-Nearest Neighbors

The MDA results for the K-Nearest Neighbors algorithm were similar to the post-MDA Bayesian results. Interestingly, the different K-variants for the *pose* and *illumination* performed the exact same. In fact, the K-Nearest Neighbor results were the same as the Bayesian results over the same dimension range. Intuitively, it is likely that this would happen when operating on a reduced data set.



5. Conclusion

Across all scenarios and data sets examined in this paper, the Bayesian classifier performed better than or equal to the K-Nearest Neighbors classifier. That is not to say that

the Bayesian classifier should be the de facto standard for image classification. The simple nature of the K-Nearest Neighbors solution makes it a practical solution when the number of samples is low and speed is a priority.

Based on the evidence present, an effective tiebreaking technique for the K-Nearest Neighbors algorithm should use the information from the previous run in order to break the tie. Of the implementations tested in this report, the `closest` method consistently performed the best, and if implemented properly, can break the tie very quickly.

As expected, Multiple Discriminant Analysis proved to be a more effective dimensionality reduction technique compared to Principal Component Analysis at lower dimensions. By preserving between-class scatter, we can project onto much lower dimensions and at the very least preserve the accuracy of the model.

Finally, it is clear that the number of training samples has a massive impact on the performance of any particular classifier, whether it be Bayesian, or non-parametric like the K-Nearest Neighbors. In circumstances where this condition holds true, dimensionality reduction techniques are less effective, and classification performance in general tends to suffer.

6. Consideration for Future Research

A more robust way to examine and rank the performance of various classifiers would have been to analyze the *top-k* error rates. In the case of the top-1 score (used in this report), one checks if the top class (the one having the highest probability, or shortest distance in the case of K-Nearest Neighbors) is the same as the target label. In the case of top-5 score, one checks if the target label is one of top 5 predictions.

The training and testing sets described in table 3 were used for all experiments conducted in this report. This, combined with the fact that the size of the training data set is small can lead to overfitting. Aside from obtaining a completely new data set with more samples, using a form of cross validation can reduce the dependence of the results on the specific training set that is used.