

Project Part 1

Pattern Recognition

ECE 759

Kudiyar Orazymbetov (korazym@ncsu.edu)
Nico Casale (ncasale@ncsu.edu)

March 16, 2018

Contents *(Note that the entries are links.)*

1	Introduction	3	3.3	Linear Discriminant Analysis (LDA)	9
			3.3.1	Algorithm steps	10
2	Feature Selection	3	3.3.2	Helpful Functions	10
2.1	Decision Tree and Extra-Tree Feature Generation	3	3.3.3	Derivation of discriminant function	11
2.2	LDA vs PCA	3			
3	Algorithm Implementations	5	4	Test Results	11
3.1	Decision Tree Algorithm	5	4.1	Preliminary Hyperparameter Optimization	11
3.1.1	Entropy and Information Gain	7	4.1.1	Decision Trees	11
3.2	Extra-Trees	9	5	Conclusion	12
			6	Code Listings	12

List of Figures

2.1	PCA implemented on 2-dimensional data [1].	4
2.2	LDA visualization [1].	4
3.1	An example decision tree.	5
3.2	Information Gain across all possible thresholds.	8

List of Tables

4.1	Lowest classification errors achieved with LDA, Decision Trees, and Extra-Trees.	11
4.2	Hyperparameter Configuration of Decision Tree for MNIST.	11

Listings

1	Main Code for Decision Tree on MNIST	12
2	Main Code for Extra-Trees on Yale B	14
3	Main Code for LDA on MNIST	16

1 Introduction

In this report, we outline the process we undertook to implement two classification algorithms: linear discriminant analysis (LDA) and Decision Trees. We tested our generated classifiers on the Yale Extended Face Dataset B and MNIST [2] datasets. To complete this, we experimented with two feature generation methods: principle components analysis (PCA) and LDA itself (for use in decision trees.) These were primarily useful in training our decision trees, which require a thorough exploration of the feature space to grow. Note that throughout this document, we use the terms '*features*' and '*attributes*' interchangeably.

After presenting an explanation of our algorithms' implementations, we will discuss our test results and preliminary exploration of the hyperparameter space. The hyperparameters are adjustments that variously control the accuracy and computational complexity of our classifiers. A more thorough exploration and validation of these parameters will be presented in the next installment of the project.

Code for this project is included in our project submission. In addition, it's available on GitHub at <https://github.com/n-casale/ece759-project>. The Git repository includes a more thorough representation of the path our code took to its current state.

2 Feature Selection

2.1 Decision Tree and Extra-Tree Feature Generation

In working with the decision trees, we utilized dimensionality reduction techniques including linear discriminant analysis (LDA) and principle components analysis (PCA) for training and testing classical decision trees, which pass over all features to consider which decision nodes to generate.

LDA reduces data with K classes in large dimensions (on the order of the number of pixels) to K dimensions which captures the most energy of the data. Implementing LDA as dimensionality reduction, MNIST and Extended YaleB datasets are reduced to 10 and 38 variables respectively (for each data point.) We then worked with this transformed data as the input to our decision tree classifier for training and testing.

In a separate experiment, we considered the effects of utilizing the raw pixel data as features in *extra-trees* which are ensembles of decision trees trained on randomly chosen features and thresholds. When testing, the mode of the predicted classes is chosen as the global prediction of the ensemble of trees. This method is shown to be more accurate than our classical decision trees trained on PCA- and LDA-generated features in Section (4).

2.2 LDA vs PCA

LDA and PCA linearly transform the data to reduce its dimensionality. Dimensionality reduction is helpful in many algorithms because it eases computational requirements and improves the generalizability of the classifier. This latter benefit is conferred by the manner in which dimensionality reduction smoothes out the individual variations in a particular class. A good dimensionality reduction algorithm should maximize the 'distance' between classes in the transformed (reduced) space, while minimizing the variance, or spread of those classes in the transformed space.

The difference between the two techniques is that LDA is supervised whereas PCA is unsupervised. PCA finds orthogonal projections which capture the most variance of each feature, thereby expressing the data along *principle components* that are the most distinct. We can visually exemplify the PCA technique as in the figure below.

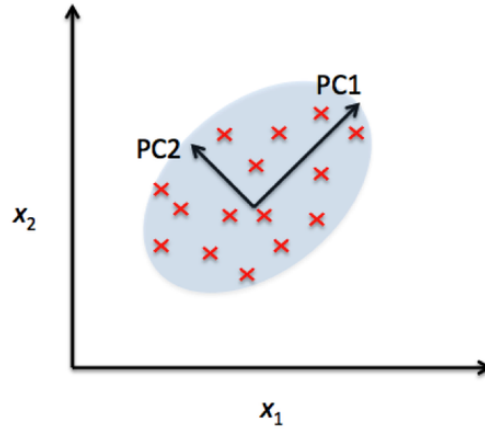


Figure 2.1: PCA implemented on 2-dimensional data [1].

In our PCA implementation, we rely on the *singular value decomposition* (SVD) to generate our features, conventionally described as *scores*. To perform the singular value decomposition and obtain our features, we must first pre-process the data. We arrange it so that the rows of our data X are the individual images in the dataset, with each column corresponding to a pixel value. Then, we subtract the column-wise empirical mean from each column to render each column zero-mean. Then, we utilize MATLAB's built-in `svd(.)` to obtain the matrices $[U, S, V]$. Our features are represented by $U * S$. Taking the transpose of this matrix, and extracting the first n columns, where n is the number of features we seek to utilize, we have a matrix of features corresponding to each image in our dataset. This method was useful, although not as effective in training decision trees as LDA, which we now discuss.

LDA maximizes the class separability and can be represented visually as in the figure below.

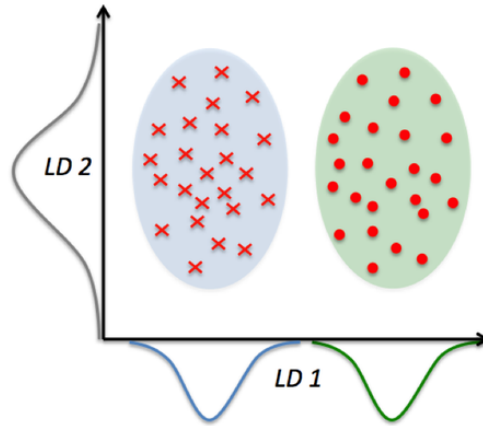


Figure 2.2: LDA visualization [1].

The underlying concept of LDA is taking the eigenvectors of $\frac{\Sigma_b}{\Sigma_w}$ where Σ_w is within-class scatter matrix and Σ_b is between-class scatter matrix. This division of matrices separates the classes away from one another. This is distinct from PCA, which is unsupervised, thereby only capturing the data in a minimum number of variables irrespective of class.

3 Algorithm Implementations

3.1 Decision Tree Algorithm

A binary decision tree is a hierarchical structure that takes input data at its root and propagates it to one of many leaves. Each *leaf* of the tree represents a class designation. To reach a leaf, the features of the data are utilized at *nodes* to make a binary decision: to proceed down the left or right *branch* of the tree? To answer this question, the node also carries a *threshold* that the feature value of the test data is compared against. If the test feature is less than the threshold, we proceed down the left branch. Otherwise, the right. An illustration of a simple decision tree is pictured below.

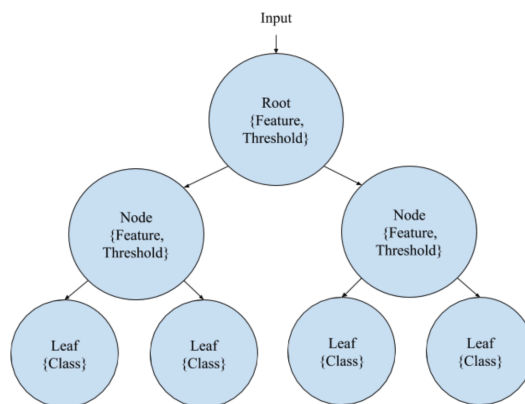


Figure 3.1: An example decision tree.

This decision tree structure needs to be generated before it can be used with test data. To train a decision tree that appropriately classifies our test data according to the features we generated, we employ a recursive function. The function signature is

```
tree = trainDecisionTree(set)
```

Where `set` is the training set, which is a MATLAB structure that contains the class labels and generated features. `tree` is the returned structure that can be used during testing. It is essentially a nested structure that contains two types of elements: nodes and leaves. At each node of the tree, a feature and threshold are specified. If a test sample's value at that particular feature is less than the threshold, the sample is passed down the left branch of the node. Similarly, if the sample's feature is greater than the threshold, it goes through the right branch. This is repeated until we reach a leaf node, which specifies a class membership.

The decision tree training algorithm has a few major steps, and proceeds by evaluating a metric called *information gain* at various configurations. For now, suffice it to say that information gain is a scalar that represents the improvement in prediction as we narrow down the set (by growing the tree) to find appropriate leaves. Our implementation most closely follows that of the ID3 and C4.5 algorithms developed by Ross Quinlan in the late 80s and early 90s [3] [4].

In general, the decision tree training algorithm is described qualitatively as follows:

1. Check stopping conditions, which generate leaves.
 - If there are no more features to split on, return a leaf with the class mode of the set.
 - The set is smaller than `minLeaf`, which is a tuning parameter that is meant to reduce overfitting of the training data. If this condition is met, return a leaf with the class mode of the set.

- If all samples in the set belong to the same class, return a leaf with the class.
 - If no feature yields an improvement to the information gain (discussed below), then return a leaf with the class mode of the set. Note that this condition is only evaluated after step 2.
2. Iterate over each feature. Sort the set along the current feature. We utilize a threshold that splits the set between adjacent feature values such that the two subsets are composed of training samples whose features are less than and greater than the threshold respectively. Because the information gain across thresholds is convex on the whole (see Fig. 3.2), we use a line search that approximates the highest information gain for each threshold. This serves to reduce the computational complexity of our training algorithm.
- Let `attributeBest` and `indexBest` be the feature and index that yield the highest information gain. Since the set is sorted, we can simply split the set at the index given by `indexBest` for the recursion.
3. Recur over the two subsets given by `indexBest` to find the next attribute that yields the highest information gain. Note that we exclude `attributeBest` in the recursion of `trainDecisionTree(.)` so that the same feature isn't chosen in the subset. In this way, the decision tree is grown so that it makes the most improvements to information gain at the nodes which are closest to the root.

The algorithm is also reproduced in pseudocode below.

Algorithm 3.1: `trainDecisionTree`

Data: *set* of training samples with class labels (*set*(1)) and attributes (features) (*set*(2)).

minLeaf, an integer specifying the minimum number of elements in *set* required to make a splitting node.

Result: *tree*, a structure containing nodes and leaves.

```

1 begin
2   check for base cases:
3   if set(2) =  $\emptyset$  then
4     no more attributes (features) to split on.
5     return leaf with mode of set(1) (class labels)
6   if length(set(1)) < minLeaf then
7     return leaf with mode of set(1) (class labels)
8   thisSetEntropy  $\leftarrow$  getEntropy(set)
9   if not thisSetEntropy then
10    all samples are in the same class.
11    return leaf with first element of set(1)
12  instantiate tracking variables:
13  attributeBest, thresholdBest, infoGainBest, indexBest  $\leftarrow$  0
14  for attribute  $\in$  range(# of features left in set) do
15    sort set along attribute.
16    perform line search heuristic to approximate max information gain by splitting set at various indices.
17    thisInfoGain  $\leftarrow$  lineSearch(set)
18    if thisInfoGain > infoGainBest then
19      attributeBest  $\leftarrow$  attribute
20      infoGainBest  $\leftarrow$  thisInfoGain
21      indexBest  $\leftarrow$  index given by line search
22      thresholdBest  $\leftarrow$  attribute value midway between indexBest and indexBest + 1
23  if not infoGainBest or not attributeBest then
24    no attribute provides an information gain.
25    return leaf with mode of set(1) (class labels)
26  re-sort set along attributeBest.
27  subsets  $\leftarrow$  getSubsets(set, attributeBest, indexBest)
28  subtree1  $\leftarrow$  trainDecisionTree(subsets(1), minLeaf)
29  subtree2  $\leftarrow$  trainDecisionTree(subsets(2), minLeaf)
30  return node with attributeBest, thresholdBest, subtree1, and subtree2

```

Once we have trained the decision tree, we can begin to test it by passing the features generated from the test set through the tree. As the features are utilized in propagating through the tree, we eventually reach a leaf node and assign its associated class to the test vector. Testing results can be found in Section (4). The testing algorithm is reproduced in pseudocode below.

Algorithm 3.2: testDecisionTree

Data: *set* of test samples with class labels (*set*(1)) and attributes (features) (*set*(2)).

tree, the struct of structs that represents the trained decision tree.

Result: *set*, the input structure modified to include predicted class membership.

```

1 begin
2   for each sample  $\in$  set do
3     treeWalked  $\leftarrow$  tree
4     classified  $\leftarrow$  false
5     attributes  $\leftarrow$  features corresponding to this sample
6     while not classified do
7       if treeWalked is a node then
8         thisAttribute  $\leftarrow$  attribute given by treeWalked
9         thisAttributeValue  $\leftarrow$  value of the sample's feature at thisAttribute
10        thisThreshold  $\leftarrow$  threshold given by treeWalked
11        remove thisAttribute from the sample's feature vector.
12        compare thisAttributeValue to thisThreshold:
13        if thisAttributeValue < thisThreshold then
14          choose left branch.
15          treeWalked  $\leftarrow$  left branch given by current treeWalked.
16        else
17          choose right branch.
18          treeWalked  $\leftarrow$  right branch given by current treeWalked.
19        else
20          treeWalked is a leaf.
21          append predicted label to set(2) at this sample.
22          classified  $\leftarrow$  true

```

With these two algorithms in place, we can train and test decision trees with varying depth and accuracy. Next, we discuss an adjacent concept that is integral to the decision tree algorithm.

3.1.1 Entropy and Information Gain

When growing the decision tree, the training algorithm utilizes a metric called *information gain* to optimize the predictive power of the tree. Before we can define information gain, we must first understand *entropy*. Entropy is an information theoretic concept that represents the amount of uncertainty in a given set of data. It is defined as

$$H(X) = - \sum_{i=1}^n P(x_i) \log_2 P(x_i) \quad (3.1)$$

Where $P(x_i)$ is the proportion of the number of elements with class x_i to the number of elements in the set X , and there are n classes in the set. Note that if all samples belong to class i , $P(x_i) = 1$ and $\log_2 P(x_i) = 0$, so the entropy is zero.

In calculating the information gain of splitting the set into two subsets, we utilize the entropy of the parent set and

subtract from it a weighted entropy of each subset [5]. This is better expressed by the following equation:

$$\text{IG}(X) = H(X) - \sum_{i=1}^2 \frac{|S_i|}{|X|} H(S_i) \quad (3.2)$$

Where $|\cdot|$ is the cardinality, or number of elements in the set, and S_i are the two subsets composed of the elements of X partitioned across a given threshold for a given feature. The figure below is a plot of the information gain across splitting indices. So, after partitioning the set across a single attribute, (all elements of S_1 are less than the threshold, and all elements in S_2 are greater), we have a scalar value of information gain to decide which splitting threshold and feature would most improve the predictive ability of the decision tree. The plot below represents just one feature, but this calculation is required across all currently available features in the set. This quickly becomes computationally taxing for a large dataset, so we utilize a line search heuristic to improve the computational time without sacrificing a significant amount of accuracy.

The line search starts in the middle of the set and splits it, computing an information gain (IG). Then, we compare this middle IG to the IGs obtained by splitting the set at the 25% and 75% indices. If either of these reveals a larger IG, we set it to the new 'middle' index and consider the IGs which emerge from splitting the set half-way between the old 'middle' and half-way between the old 'left' or 'right' (depending on which one yielded a larger IG.) This allows us to approximate the maximum value of the information gain for a given feature without a brute-force technique of considering each threshold. We ensured that the line search incorporates a variety of locations to coax it to approaching the global maximum, rather than getting caught in the minor variations between adjacent indices. Please consult `code/decisionTree/trainDecisionTree.m` for further details.

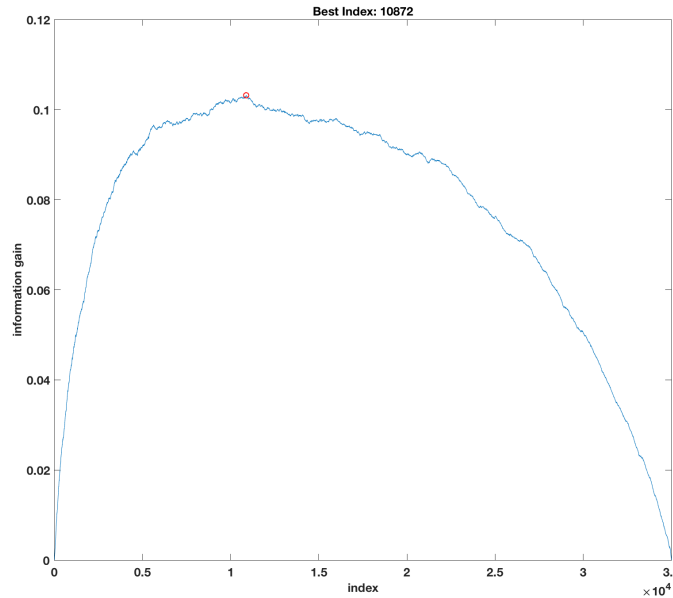


Figure 3.2: Information Gain across all possible thresholds.

3.2 Extra-Trees

The relative inadequacy of decision trees in classifying the samples of MNIST and Yale B led us to explore other tree-like classifiers. The first alternative that we considered is called *extra-trees*, which are ensembles of binary decision trees that are grown in a stochastic manner. In testing, these ensembles *vote* on the class by propagating the test sample down each tree until a leaf is found. The mode of all the votes is assigned to the global class prediction. See Section (4) for the result of this experiment.

Training extra-trees requires less computational and theoretical effort than classical decision trees. The reduced steps are certainly tangential to a basic decision tree, but incorporate less rigour. By utilizing the ensemble of trees, we can overcome the inaccuracies of a single random decision tree and make a strong prediction.

The steps to train an extra-tree are as follows:

1. Check stopping conditions, which generate leaves.
 - If there are no more features to split on, return a leaf with the class mode of the set.
 - The set is smaller than `minLeaf`, which is a tuning parameter that is meant to reduce overfitting of the training data. If this condition is met, return a leaf with the class mode of the set.
 - If all samples in the set belong to the same class, return a leaf with the class.
2. Choose a random feature. In extra-trees, we don't need to generate features, we can simply use the raw pixels as features.
3. Find the mean and variance of this feature across all samples in the set. Generate a random value from a normal distribution with this mean and variance.
4. Recur these steps on the subsets obtained by splitting the parent set on the randomly chosen feature and threshold, where the first subset contains samples whose feature is less than the threshold, and the second subset contains those which are greater than the threshold.

For MNIST, we generated 100 extra-trees, trained on half of the dataset. In testing, the majority-vote of the trees was used to assign a predicted class to the test vectors. Likewise, for Yale B, we utilized half the dataset and 100 extra-trees.

3.3 Linear Discriminant Analysis (LDA)

Our classification success criterion is the extent to which the error rate is minimized. For LDA, we employ the Bayes Rule to classify our test points after training the classifier. We assign the test point to the class with the highest conditional probability (*i.e.* $P(w_i|x)$, where w_i is the class and x the test vector). In practice, it is not feasible to get conditional probability for a given point unless we have a large amount of data. So we estimate the distribution and calculate the probabilities from there.

LDA relies on the assumption of a normal distribution for each class. Linear discriminant analysis generally achieves good performance in the tasks of face and object recognition, even though the assumptions of common covariance matrix among groups and normality are often violated [6].

Since the MNIST and Yale datasets are high-dimensional, we cannot check the normality of individual pixels. Instead, we reduce the dimensions via orthogonal projection. To do this, we employ multiclass LDA, as we have several classes (10 for MNIST and 38 for Yale B.) The class separation in a direction \mathbf{w} in this case is given by the ratio [7].

$$S = \frac{\mathbf{w}^T \Sigma_b \mathbf{w}}{\mathbf{w}^T \Sigma \mathbf{w}} \quad (3.3)$$

In the case of two classes, this reduces to the ratio of between-class variance and within-class variance.

3.3.1 Algorithm steps

Below, we describe the steps necessary to construct the LDA classifier. We begin by instantiating some variables with notation:

n is the number of classes.

N is the total number of training data samples.

N_i is the number of points in each class.

μ_i and μ are mean vectors for each class and global mean for the data.

The steps that we follow in our LDA algorithm are as follows.

1. We calculate within-class scatter matrix for each class

$$\Sigma_i = \frac{1}{N_i - 1} \sum_{\mathbf{x} \in D_i} (\mathbf{x} - \mu_i) (\mathbf{x} - \mu_i)^T \quad (3.4)$$

then sum them to obtain

$$\Sigma_W = \sum_{i=1}^n (N_i - 1) \Sigma_i \quad (3.5)$$

2. We then find the average within-class scatter matrix by calculating

$$\Sigma = \frac{\Sigma_W}{N} \quad (3.6)$$

3. We also calculate the between-class scatter matrix by

$$\Sigma_B = \sum_{i=1}^n \frac{N_i}{N} (\mu_i - \mu) (\mu_i - \mu)^T \quad (3.7)$$

4. We need to find eigenvectors and eigenvalues of $\Sigma^{-1}\Sigma_B$.
5. We then sort the eigenvectors depending on the magnitude of eigenvalues.
6. The number of highest eigenvalues will be $c - 1$ which will be 9 and 37 respectively for MNIST and Extended Yale datasets.
7. Project our data onto the subspace(constructed by the eigenvectors of the highest eigenvalues).
8. Since we have reduced dimensions of our data, we can easily apply a discriminant function for a test vector to see which class it belongs
9. The discriminant function for each class is

$$f_i(x_k) = \mu_i^T w_a^{-1} x_k^T - \frac{1}{2} \mu_i^T w_a^{-1} \mu_i^T + \ln(P_i) \quad (3.8)$$

where P_i is the probability of each class.

10. Then we select the class with highest discriminant function evaluation.

3.3.2 Helpful Functions

In the above pseudo code, we used own created functions. After we do steps 1-7, there is a function `classify_comparison` which applies discriminant function in order to classify the test data. I also tried to calculate if I get good results using the distances from the centers of each classes using `classify_from_centroid`, although it did not work very well on MNIST dataset it does perform well with Extended Yale datasets.

3.3.3 Derivation of discriminant function

We classify the data point as being in class i if $P(\mathbf{x}|i)P(i) > P(\mathbf{x}|j)P(j)$, for $\forall j \neq i$. Assuming all covariances are equal $\Sigma = \Sigma_i = \Sigma_j$ and they follow a Multivariate Normal distribution formula $P(\mathbf{x}|i) = \frac{1}{(2\pi)^{n/2}|\Sigma_i|^{1/2}} \exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i))$ the above condition becomes $\frac{P(i)}{(2\pi)^{n/2}|\Sigma_i|^{1/2}} \exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)) > \frac{P(j)}{(2\pi)^{n/2}|\Sigma_j|^{1/2}} \exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j))$ Taking the logarithm of both sides we get $\ln(|\Sigma|) - 2\ln(P(i)) + (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) < \ln(|\Sigma|) - 2\ln(P(j)) + (\mathbf{x} - \boldsymbol{\mu}_j)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_j)$ After some algebraic manipulation, we have $\ln(P(i)) + \boldsymbol{\mu}_i \Sigma^{-1} \mathbf{x}^T - \frac{1}{2} \boldsymbol{\mu}_i \Sigma^{-1} \boldsymbol{\mu}_i^T > \ln(P(j)) + \boldsymbol{\mu}_j \Sigma^{-1} \mathbf{x}^T - \frac{1}{2} \boldsymbol{\mu}_j \Sigma^{-1} \boldsymbol{\mu}_j^T$ Let $f_i = \boldsymbol{\mu}_i \Sigma^{-1} \mathbf{x}^T - \frac{1}{2} \boldsymbol{\mu}_i \Sigma^{-1} \boldsymbol{\mu}_i^T + \ln(P(i))$, then we can write $f_i > f_j, \forall j \neq i$

4 Test Results

The table below illustrates the lowest classification errors we were able to achieve with our Decision Tree and LDA classifiers on the MNIST and Yale Extended B datasets. In LDA, we split the datasets in half so we can the other half for test. Using decision trees on MNIST, we utilized half of the dataset for training and half for testing. However, since the Yale B dataset is so small, we were only able to get lower error rates by using $\sim 95\%$ of the set for training and $\sim 5\%$ for testing (on decision trees.) For extra-trees, we split the two datasets in half to achieve our results.

Algorithm	MNIST	Yale B
LDA	13.5%	6.8%
Decision Tree	17.42%	57.89%
Extra-Trees	4.89% (100 extra-trees)	34.96% (100 extra-trees)

Table 4.1: Lowest classification errors achieved with LDA, Decision Trees, and Extra-Trees.

4.1 Preliminary Hyperparameter Optimization

4.1.1 Decision Trees

The table below illustrates some results that we obtained during various configurations of the decision tree training algorithm on MNIST. In an effort to optimize training time vs. test accuracy, we noticed some trends in the algorithm's performance. Notably, utilizing `minLeaf` doesn't seem to contribute to over-fitting, as was noted in other implementations [8]. This may be in part to the type of algorithm we are using. Particularly, the pure C4.5 algorithm moves back up the tree if certain conditions are met, thereby improving a previously-made split. This bore more complexity than we were able to achieve at this time, but may be responsible for a more simpler tree generation on our part. Note that for MNIST, we used half of the 70k samples for training and half for testing.

Training Configuration	1	2	2	2	3	3	1	1	1	1	4	4
numFeatures	30	30	60	200	200	200	30	100	100	20	10	10
minLeaf	1	1	2	2	3	4	4	4	1	1	1	3
Error Rate (%)	22	26	26	26	23	24	22.6	23.6	23.5	22.9	17.42	20.37
Mins To Train	48	2	2	3	12	12	47	60	102	11	3	3

Table 4.2: Hyperparameter Configuration of Decision Tree for MNIST.

The training configurations were primarily modifications to the algorithm's search through the features to find the best one to split the set on. They are as follows:

1. All features were considered for the best information gain in a particular set.
2. The decision tree was grown by considering only the first feature in the set. This was an attempt to reduce the computational complexity.
3. If the number of features left in a set was greater than 5, we considered only the first 5. Otherwise, we considered all features left.
4. For configurations 1-3, we utilized features generated via PCA. For this experiment, we utilized features generated with LDA. This yielded our best performance.

5 Conclusion

In conclusion, we have shown that utilizing Linear Discriminant Analysis, Decision Trees, and Extra-Trees on the Extended Yale Dataset B and MNIST datasets yields a favorable result. In the next part of this project, we will improve our classification accuracy and partake in a more thorough exploration of the hyperparameter space to optimize and generalize our classifiers. Additionally, we will develop a stand-alone demonstration of our algorithms' process from training to testing.

References

- [1] Sebastian Raschka. What is the difference between lda and pca for dimensionality reduction?, 2018. [Online; accessed 2018-03-14].
- [2] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010.
- [3] Wikipedia contributors. Id3 algorithm — wikipedia, the free encyclopedia, 2017. [Online; accessed 14-March-2018].
- [4] Wikipedia contributors. C4.5 algorithm — wikipedia, the free encyclopedia, 2018. [Online; accessed 14-March-2018].
- [5] Wikipedia contributors. Information gain in decision trees — wikipedia, the free encyclopedia, 2017. [Online; accessed 15-March-2018].
- [6] Tao Li, Shenghuo Zhu, and Mitsunori Ogihara. Using discriminant analysis for multi-class classification: an experimental investigation. *Knowledge and Information Systems*, 10(4):453–472, Nov 2006.
- [7] Wikipedia contributors. Linear discriminant analysis — wikipedia, the free encyclopedia, 2017. [Online; accessed 15-March-2018].
- [8] Mathworks. Matlab documentation - fitctree, 2018. [Online; accessed 2018-03-14].

6 Code Listings

Below are the primary scripts that execute the project. Please see the `code` folder for supporting functions.

Listing 1: Main Code for Decision Tree on MNIST

```
1 %{\n2
```

```
3 kudiyar orazymbetov
4 n casale
5
6 ECE 759 Project
7 18/03/16
8
9 this script orchestrates the training and testing of
10 the decision tree classifier
11
12 features are generated using principal component analysis
13
14 %}
15
16 clear; close all;
17 addpath('utility', 'MNIST', 'MNIST/data', 'MNIST/loadMNIST', 'lda');
18
19 fprintf('begin MNIST decision tree script\n');
20
21 % hyper-parameters
22 N_tr = 35e3; % training samples
23 N_te = 35e3; % test samples
24
25 % for feature selection
26 usePCA = false;
27 if usePCA
28     numFeatures = 30;
29 else % use LDA
30     numFeatures = 10;
31 end
32
33 % for decision tree
34 minLeaf = 1; % to prevent overfitting
35
36 %seed = 152039828;
37 %rng(seed); % for reproducibility
38
39 % partition data
40 % MNIST contains 70k examples
41 [train, test] = loadMNIST(N_tr);
42
43 %% dimensionality reduction / feature generation
44 % via prinicpal component analysis (pca) (svd)
45 st = cputime;
46
47 if usePCA
48     [train, U, V] = pca_(train, numFeatures);
49 else
50     [train, test] = lda_features(train, test, 0:numFeatures-1);
51 end
52
53 fprintf('Features Generated in %4.2f minutes\n', (cputime - st)/60);
54
55 %% train
56 st = cputime;
57
58 tree = trainDecisionTree({train{2:3}}, minLeaf);
```

```

59
60 fprintf('Trained in %4.2f minutes\n',(cputime - st)/60);
61
62 %% test
63 st = cputime;
64
65 if usePCA
66     test{3} = (U'*test{1}*'V)';
67     test{3} = test{3}(1:numFeatures,:);
68 end
69
70 test = testDecisionTree(test, tree);
71
72 fprintf('Tested in %4.2f minutes\n', (cputime - st)/60);
73
74 % Classification Error
75 errors = nnz(test{2}(:,1) ~= test{2}(:,2));
76 errorRate = (errors/N_te)*100;
77
78 filename = sprintf('mnist_tree%2.0f_%d.mat', errorRate, minLeaf);
79 %save(filename, 'tree');
80 %filename = sprintf('UV%2.0f.mat', errorRate);
81 %save(filename, 'U', 'V', '-v7.3');
82
83 fprintf('numFeatures: %d, minLeaf: %d, error rate: %2.2f\n', ...
84         numFeatures, minLeaf, errorRate);

```

Listing 2: Main Code for Extra-Trees on Yale B

```

1  %{
2
3  kudiyar orazymbetov
4  n casale
5
6  ECE 759 Project
7  18/03/16
8
9  this script orchestrates the training and testing of
10 the extra-tree classifier
11
12 there are no features used besides the raw pixels
13
14 %}
15
16 clear; close all;
17 addpath('utility', 'YaleB', 'YaleB/data');
18
19 fprintf('begin Yale B extra tree script\n');
20
21 % hyper-parameters
22 N_tr = 2000; % training samples
23 N_te = 2414 - N_tr; % test samples
24
25 % for extra tree
26 minLeaf = 5; % to prevent overfitting
27 numTrees = 100; % ensemble for majority voting
28

```

```

29 %seed = 152039828;
30 %rng(seed); % for reproducibility
31
32 % partition data
33 % MNIST contains 70k examples
34 [train, test] = loadYaleB(N_tr);
35
36 % features are the raw pixels, so we reorder the cell
37 train = {train{2}, train{1}};
38 test = {test{2}, test{1}};
39
40 %% train
41 st = cputime;
42
43 % create an ensemble of random trees
44 trees = cell(numTrees, 1);
45 for tree = 1:numTrees
46
47     fprintf('tree: %d\n', tree);
48     trees{tree} = trainExtraTree(train, minLeaf);
49
50 end
51
52 fprintf('Trained %d extra-trees in %4.2f minutes\n', numTrees, (cputime - st)
    /60);
53
54 %% test
55 st = cputime;
56
57 test = testExtraTrees(test, trees);
58
59 fprintf('Tested in %4.2f minutes\n', (cputime - st)/60);
60
61 % Classification Error
62 errors = nnz(test{1}(:,1) ~= test{1}(:,2));
63 errorRate = (errors/N_te)*100;
64
65 filename = sprintf('mnist_extratree%2.0f_%d.mat', errorRate, minLeaf);
66 %save(filename, 'trees');
67
68 fprintf('minLeaf: %d, error rate: %2.2f\n', ...
69     minLeaf, errorRate);
70
71 %% Test Results
72 %{
73
74     , * * * & & , , , ,
75 numFeatures : 30 30 60 200 200 200 30 100 100 20 10 (lda)
76 minLeaf : 1 1 2 2 3 4 4 4 1 1 1
77 Error Rate : 22 26 26 26 23 24 22.6 23.6 23.5 22.9 18.42
78 Mins To Train : 48 2 2 3 12 12 47 60 102 11 3
79
80 , - all attributes were considered for the best information gain for a
81 particular set.
82
83 * - for these experiments, the decision tree was grown by considering
84 only the first attribute in the set. all things held equal, this was simply

```

```

84 an attempt to reduce the computational complexity.
85
86 & - if the number of attributes was > 5, consider only the first 5, else,
87 consider all attributes (1-5)
88
89 %}

```

Listing 3: Main Code for LDA on MNIST

```

1  %{
2
3  kudiyar orazymbetov
4  n casale
5
6  ECE 759 Project
7  18/03/16
8
9  this script orchestrates the training and testing of
10 the linear discriminant analysis classifier
11
12 %}
13
14 clear; close all;
15 addpath(' ../utility', ' ../MNIST', ' ../MNIST/data', ' ../MNIST/loadMNIST', ...
16         ' ../lda');
17
18 seed = 152039828;
19 rng(seed); % for reproducibility
20
21 % define parameters
22 N_tr = 35e3; % training samples
23 N_te = 35e3; % test samples
24 k = 10; % number of classes
25 % partition data
26
27 %{
28     1/2 of the dataset should be for training
29     the other for testing
30
31     MNIST contains 70k examples
32 %}
33
34 [train, test] = loadMNIST(N_tr);
35 %% use PCA to see the results
36 % numFeatures = 20;
37 % [train, U, V] = pca_(train, numFeatures);
38
39 %% Construct scatter matrices and calculate within-class and between class
40 % covariance
41 mu = mean(train{1,1}, 2);
42 num_variables = size(train{1,1},1);
43 % Let's standardize the data;
44 %variance = var(train{1,1}, 0,2);
45 %train{1,1} = (train{1,1}-repmat(mu,1, 60000))./variance;
46
47 %
48 Si = zeros(num_variables); Sb = zeros(num_variables);

```



```

49 S_cov = zeros(num_variables);
50 for i = 0:k-1
51     ind = (train{1,2} == i);
52     N_i = sum(ind);
53     x = train{1,1}(:, ind);
54     mu_i = mean(x, 2);
55     S_cov = S_cov + cov(x');
56     Si = Si + (1/N_tr)*(x - (repmat(mu_i,1, N_i)))*(x - (repmat(mu_i,1, N_i)))
        ' ;
57     Sb = Sb + (N_i/N_tr)*(mu_i - mu)*(mu_i - mu)'; % (1/k)
58 end
59
60 % We apply singular value decomposition in order to find eigenvalues and
61 % eigenvectors
62 [U D V] = svd(pinv(Si)*Sb); % lets try S_cov/k instead of Si; but it is the
    same result
63 a = [];
64 for i = 1:(k)
65     a = [a D(i,i)];
66 end
67
68 % from here we can see that we only have 9 highest values as we expected
69
70 % We transform the training and testing data to a subspace
71 transf_matrix = U(:,1:(k-1));
72 transf_train = train{1,1}'* transf_matrix;
73 transf_test = test{1,1}'*transf_matrix;
74 % We find mean vector and covariance matrix for each class
75 mu_each_class = zeros(k-1, k);
76 cov_each_class = {};
77 sum_cov = zeros(k-1,k-1);
78 for i =0:k-1
79     ind = find(train{1,2} == i);
80     X = transf_train(ind, :);
81     mu_each_class(:, i+1) = mean(X, 1)';
82     cov_each_class{1, i+1} = cov(X);
83     sum_cov = sum_cov + cov(X);
84 end
85 % since we assume equal covariance in all classes, we take the average of
86 % covariance matrices
87 average_cov = sum_cov/k;
88 cov_equal_each_class = {average_cov average_cov average_cov average_cov
    average_cov average_cov average_cov average_cov average_cov average_cov};
89 % this part is just a test on how nearest neighbors work
90 % parfor n = 1:13
91 % % we apply Nearest neighbors in order to find which class it belongs
92 %     accuracy(n) = classifyNN(n,transf_test', transf_train', test{1,2}, train
    {1,2});
93 % end
94 [acc_test_comp acc_train_comp] = classify_comparison(k,5,mu_each_class,
    cov_equal_each_class, transf_test', test{1,2}, transf_train', train{1,2});
    % 0.88 and 0.89 resp using just knn
95
96 %% This part uses kNN in order to class after transformation
97 [acc_test_5 acc_train_5] = classifyNN(k,5,mu_each_class, cov_each_class,
    transf_test', test{1,2}, transf_train', train{1,2}); % 0.88 and 0.89 resp

```

```

    using just knn
98 [acc_test_5 acc_train_5] = classifyNN(k,5,mu_each_class, cov_equal_each_class,
    transf_test', test{1,2}, transf_train', train{1,2}); % 0.87 and 0.88 resp
    using just knn
99 [acc_test_5_p acc_train_5_p] = classifyNN_pure(5,transf_test', transf_train',
    test{1,2}, train{1,2});
100
101 % cov each class separately works better in each case
102 % we plot the results to see the best number of nearest neighbors
103 % [acc_test_5 acc_train_5] = classifyNN(5,transf_test', transf_train', test
    {1,2}, train{1,2}); % 0.86 and 0.85 resp using just knn
104 % this is my first attempt
105
106 f = instantiateFig(1);
107 plot([1:13],accuracy*100, 'r.')
108 prettyPictureFig(f);
109 xlabel('Nearest neighbor number');
110 ylabel('Accuracy of test model');
111
112 print('.../..images/NN after LDA', '-dpng');
113 % instead we can use Euclidean distance metric to evaluate the classes by
114 % calculating the distances from each class centroid
115 centroid = zeros(k, k-1);
116 for i = 0:k-1
117     ind = (train{1,2} == i);
118     N_i = sum(ind);
119     centroid(i+1, :) = mean(transf_train(ind,:), 1);
120 end
121
122 accuracy1_train = classify_from_centroid(transf_train', train{1,2}, centroid);
123 accuracy1 = classify_from_centroid(transf_test', test{1,2},centroid);
124
125 %% this is classification through matlab discriminant classification
126 mdl = fitcdiscr(transf_train, train{1,2},'DiscrimType','linear'); % this one
    works since n>m in tansf_train
127 mdl = fitcdiscr(train{1,1}', train{1,2},'DiscrimType','linear'); % error saying
    Predictor x1 has zero within-class variance.
128 pred = predict(mdl, transf_test);
129 count = 0;
130 for i = 1:size(transf_test,1)
131     if pred(i) == test{1,2}(i);
132         count = count + 1;
133     end
134 end
135 acc = count/size(transf_test,1) % 0.8639
136 % we get a similar result as our methodp

```