

# Project Part 1

## Pattern Recognition

### ECE 759

Kudiyar Orazymbetov ([korazym@ncsu.edu](mailto:korazym@ncsu.edu))  
Nico Casale ([ncasale@ncsu.edu](mailto:ncasale@ncsu.edu))

March 13, 2018

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

---

<b>1</b>	<b>Introduction</b>	<b>2</b>	<b>3.1</b>	Decision Tree Algorithm . . . . .	<b>2</b>		
			<b>3.1.1</b>	Information Gain and Entropy . .	<b>3</b>		
<b>2</b>	<b>Feature Selection</b>	<b>2</b>					
	2.1	Decision Tree Feature Generation . . . . .	<b>2</b>	<b>4</b>	<b>Algorithm Implementations</b>	<b>3</b>	
		2.1.1	MNIST . . . . .	<b>2</b>	<b>4.1</b>	Linear Discriminant Analysis . . . . .	<b>4</b>
<b>3</b>	<b>Algorithm Implementations</b>	<b>2</b>	<b>5</b>	<b>Code Listings</b>	<b>4</b>		

## List of Figures

---

## Listings

---

# 1 Introduction

---

## 2 Feature Selection

---

### 2.1 Decision Tree Feature Generation

---

#### 2.1.1 MNIST

In working with the decision trees, we utilized the SVD of each image in in the training set.

## 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 attribute of the test data is compared against. If the test attribute is less than the threshold, we proceed down the left branch. Otherwise, the right.

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 raw data (unused), 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, an attribute and threshold are specified. If a test sample's feature at that particular attribute 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 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.

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. Because the information gain across thresholds is convex on the whole (see Fig. 3.1), we use a line search that approximates the highest information gain for each threshold.

Let `attributeBest` and `indBest` 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 `indBest` for the recursion.

3. Recur over the subsets given by `indBest` to find the next attribute that yields the highest information gain. Note that we exclude the attribute we chose in this execution of `trainDecisionTree(.)`.

This can be expressed in psuedocode as

---

**Algorithm 3.1:** Train Decision Tree

---

**Data:** *set* of training samples with class labels (*set2*) and features (*set1*).

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

---

```

1 begin
2    $V \leftarrow U$ 
3    $S \leftarrow \emptyset$ 
4   for  $x \in X$  do
5      $NbSuccInS(x) \leftarrow 0$ 
6      $NbPredInMin(x) \leftarrow 0$ 
7      $NbPredNotInMin(x) \leftarrow |ImPred(x)|$ 
8   for  $x \in X$  do
9     if  $NbPredInMin(x) = 0$  and  $NbPredNotInMin(x) = 0$  then
10       $AppendToMin(x)$ 
11  while  $S \neq \emptyset$  do
12    remove  $x$  from the list of  $T$  of maximal index
13    while  $|S \cap ImSucc(x)| \neq |S|$  do
14      for  $y \in S - ImSucc(x)$  do
15        { remove from  $V$  all the arcs  $zy : \}$ 
16        for  $z \in ImPred(y) \cap Min$  do
17          remove the arc  $zy$  from  $V$ 
18           $NbSuccInS(z) \leftarrow NbSuccInS(z) - 1$ 
19          move  $z$  in  $T$  to the list preceding its present list
20          {i.e. If  $z \in T[k]$ , move  $z$  from  $T[k]$  to  $T[k - 1]$ }
21         $NbPredInMin(y) \leftarrow 0$ 
22         $NbPredNotInMin(y) \leftarrow 0$ 
23         $S \leftarrow S - \{y\}$ 
24         $AppendToMin(y)$ 
25       $RemoveFromMin(x)$ 

```

---

### 3.1.1 Information Gain and Entropy

Figure 3.1: Information Gain across all possible thresholds.

## 4 Algorithm Implementations

---

## 4.1 Linear Discriminant Analysis

---

Our classification criterion is to misclassify as small as possible. The rule we employ in classifying the data points is through use of Bayes Rule. We put the data point to the group with the highest conditional probability. In practice, it is not feasible to get conditional probability for a given point unless we have a huge data. So we should assume the distribution and calculate the probabilities.

LDA relies on the assumption of normal distribution of data for each class. Linear discriminant analysis frequently achieves good performances in the tasks of face and object recognition, even though the assumptions of common covariance matrix among groups and normality are often violated (Duda, et al., 2001) (Tao Li, et al., 2006)'.

Since MNIST and Yale datasets are high-dimensional, we can not check the normality of variables. Instead, we can reduce the dimensions via a projection. We need to employ multiclass LDA as we have 10 classes. The class separation in this case will be given by the ratio of  $\frac{\mathbf{w}^T \Sigma_b \mathbf{w}}{\mathbf{w}^T \Sigma_w \mathbf{w}}$ . In the case of two classes, this will reduce just to the ratio of between-class variance to within-class variance.

The steps that we follow in our LDA algorithm are:

1. We calculate within-class scatter matrix  $\Sigma_i = \frac{1}{N_i-1} \sum_{\mathbf{x} \in D_i}^n (\mathbf{x} - \boldsymbol{\mu}_i) (\mathbf{x} - \boldsymbol{\mu}_i)^T$  for each class, then sum them up to get  $\Sigma_W = \sum_{i=1}^c (N_i - 1) \Sigma_i$ .
2. We also calculate the between-class scatter matrix by  $S_B = \sum_{i=1}^c N_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$ .
3. We need to find eigenvectors and eigenvalues of  $\Sigma^{-1} \Sigma_b$ .
4. We then sort the eigenvectors depending on the magnitude of eigenvalues.
5. The number of linear discriminants will be  $c - 1$  which will be 9. We need to verify that.
6. Project our data onto the subspace(constructed by the eigenvectors of the highest eigenvalues).
7. Since we have a reduced dimension for our data, we can easily apply nearest neighbors method in order to get the classes for our test data
- 8.

## 5 Code Listings

---

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