

# Assignment 2: Data Mining and Machine Learning

2ID90/2IDC0 Q3, 2013-2014

March 14, 2014

## 1 The Data

In this assignment you are going to build classifiers for recognizing handwritten digits. These digits are taken from the online MNIST database at <http://yann.lecun.com/exdb/mnist>; some example digits are shown in Figure 1.

The database consists of 60.000 training patterns and 10.000 test patterns. More information on the data can be found at the given URL. Sample java code for reading the datasets as well as the data itself will be provided.

## 2 K-nearest neighbour Classifier

1. As a first classifier you are asked to construct a k-NN classifier. To do so you are asked to implement the Classifier interface, as specified below. This interface is part of a java package<sup>1</sup> that will be provided. Initially, the getConfusionMatrix method may simply return null.

```
1 public interface Classifier<V extends Features, Target> {
2     /** returns the classification of feature vector v. */
3     Target classify(V v);
4
5     /**
6      * @ return errorRate of this classifier for the testData.
7      */
8     double errorRate(Dataset<V, Target> testData);
9
10    /** returns a map m such that m.get(a).get(b) is the
11     * number of times target a has been recognized as b.
12     * @param testData dataset for which matrix is computed
13     * @return confusion matrix
14     */
15    Map<Target, Map<Target, Integer>>
16    getConfusionMatrix(Dataset<V, Target> testData);
17 }
```

2. Run the k-NN classifier with the given training and test patterns<sup>2</sup>. Note down the resulting errorRate for different values of k. Analyze the behaviour of this classifier by looking at the confusion matrix.<sup>3</sup>

## 3 Decision tree classifier

1. Next you are going to build a decision tree classifier.
  - (a) As a preparation step implement and test the following class.



Figure 1: some labeled hand-written digits in the MNIST database.

<sup>1</sup> Class KNN in package nl.tue.s2id90.classification.knn has an empty implementation of this interface. You may adapt this class to fit your needs.

<sup>2</sup> See section 4.1 for instructions on how to read the MNIST database.

<sup>3</sup> If you implement the getConfusionMatrix method of Classifier, you can visualize the confusion matrix using method showIt of ConfusionMatrixPanel.

```

1 public class Information {
2     /**
3      * @param p probability distribution
4      * @return the entropy of p.
5      */
6     public static double entropy(double ... p){
7         // YOUR IMPLEMENTATION
8     }
9     /**
10    * @param x argument
11    * @return 2log(x)
12    */
13    public static double log2(double x){
14        // YOUR IMPLEMENTATION
15    }
16    /**
17    * @param x argument
18    * @return x 2log(x)
19    */
20    public static double xlog2x(double x) {
21        // YOUR IMPLEMENTATION
22    }
23 }
    
```

(b) The next step is to write a program to build a decision tree for discrete variables only. See section?? for a description of helpful provided software. Eventually,

- i. test on toy data<sup>4</sup>;
- ii. test on digit database; discuss your results.

<sup>4</sup> There are two toy datasets; see section 4.2 on how to access them.

(c) Implement pruning of a decision tree and apply to the MNIST tree.

## 2. Build decision tree with continuous variables

(a) decide on granularity of data (speed vs accuracy?)

3. Build several decision trees from (randomly?) chosen subsets of the training data. Combine their classification results to build a boosted classifier. Optimize your results by choosing proper parameter settings.

4. Compute additional features to further improve the error rate.

## 4 Software

### 4.1 Read MNIST database

For reading the MNIST database (static) methods in the class `HandWrittenDigits` are used. The data read from the database consists of labeled images, which have image data encoded in unsigned bytes and a classification label encoded in a byte.

```

1 List<LabeledImage> trainingImages , testImages ;
  //read 15000 training records and shuffle them
3 trainingImages = HandWrittenDigits.getTrainingData(15000,true);

5 //read test data
  testImages = HandWrittenDigits.getTestData();
    
```

The resulting lists of images may be viewed using an auxiliary method in DigitsUtil. The next line of code shows the first 25 test images in a 5 column grid of images; see Figure 2.

```

DigitsUtil.showImages("test images",testImages.subList(0,25),5)
    
```

For usage in the Classifier interface this data has to be reordered in a Map that maps features (here of type ImageFeatures<Double>) on labels (here of type Byte).

```

1 Map<ImageFeatures<Double>, Byte> trainingDataset , testDataset ;
  trainingDataset= new HashMap<>();
3 for (LabeledImage image : trainingImages) {
    trainingDataset.put(new Doubles(image), image.getLabel());
5 }

7 testDataset = new HashMap<>();
  for (LabeledImage image : testImages) {
9     testData.put(new Doubles(image), image.getLabel());
  }
    
```



Figure 2: The first 25 test images in the MNIST database.

For a subtype F of Features and a label type L, a typical usage of a labeled dataset is given below.

```

Map<F extends Features, L> dataset;
2 ...
  for (Entry<F,L> train : dataset.entrySet()) {
4     F f = train.getKey();
    L label = train.getValue()
6     ...
  }

8 // or alternatively
10 for (F train : dataset.keySet()) {
    L label = trainingDataset.get(F);
12 ...
  }
    
```

Computing the distance between two ImageFeatures might be done as follows.

```

1 public double distance(
  ImageFeatures<Double> fo , ImageFeatures<Double> f1
3 ){
  Double[] ao = fo.getValues(), a1 = f1.getValues();
5   final int n = fo.getNumberOfAttributes(); double d = 0;
  for (int i = 0; i < n; i = i + 1)
7       double h = ao[i] - a1[i]; d += h * h;
  return d;
9 }
    
```

Alternatively, you may use `fo.get(i)` to get the value of the *i*-th feature of `fo`.

#### 4.2 Toy data

There are two toy datasets: `SkiData` and `GolfData`. They can be used as in the following code fragment; the resulting tree is shown in Figure 3.

```

1 // create a proper dataset for decisiontree building
2 public static LabeledDataset2<Features , GolfData.PLAY>
3   getDataset( GolfData data )
4 {
5   LabeledDataset2<Features , GolfData.PLAY> set ;
6   set = new MyLabeledDataset2() ;
7   set.putAll( data.getClassification() ) ;
8   return set ;
9 }
10
11 // build decision tree and show it on screen
12 public static void main( String [] argv ) {
13   GolfData d = new GolfData() ;
14   LabeledDataset2<Features , GolfData.PLAY> dataset = getDataset( d ) ;
15   DecisionTree<Features , GolfData.PLAY> dt ;
16   dt = new MyDecisionTree<>( dataset ) ;
17   DotUtil.showDotInFrame( dt.toDot() , "golf data" ) ;
18 }
    
```

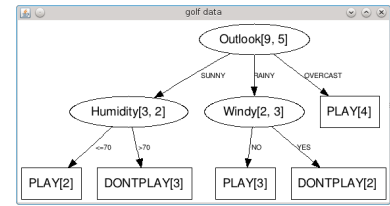


Figure 3: Golf data decision tree created using an online GraphViz server. Note that large trees are not accepted by this server.

#### 4.3 class *LabeledDataset*

In order to program your own decision tree builder you are given several classes to give you a quick start. In this section and the following sections the most important ones are discussed.

The class `LabeledDataset` is a container for feature vectors and their labels. It is fully implemented. Some of its methods are:

- `getLabel(f)` that returns the label of feature vector *f*;
- `getFeatureVectors()` that returns a list of all feature vectors in the dataset;
- `getLabels()` that returns a set of all labels in the dataset.

#### 4.4 class *LabeledDataset2*

For decision tree building it is handy to have a mapping from a label to feature vectors with that label. This functionality is supplied by `LabeledDataset2` in methods like:

- `getFeatures(L label)` that returns all feature vectors with classification label.

The class `LabeledDataset2` also contains several methods that you have to implement for building decision trees; for example:

```

2  /** splits dataset in subsets with a constant value
    * for the i-th feature.
    * @return result of splitting, the keys in this map are
    * the values of the i-th attribute.
    */
6  public Map<Object, LabeledDataset2<V, L>> discreteSplit(int i)
    { ... }
8
10 /** returns the probabilities of the classification labels
    * for this dataset.
    * @return probability of the labels.
    */
12 public double[] classProbabilities() { ... }
14
16 /**
    * @param index index of attribute used for discrete splitting.
    * @return gain by discrete splitting on attribute index.
    */
18 public double gain(int index) { ... }
20 ...
    
```

#### 4.5 class DecisionTree

A decision tree is a tree with labels, corresponding to attribute values, on the edges<sup>5</sup>. For this functionality class LabeledTree has been implemented. It contains functionality, among others, for inspecting the labels and walking the tree. The class DecisionTree inherits from LabeledTree and uses Objects as edge labels (that is, any possible attribute value). Furthermore, each node in the tree contains a labeled dataset with feature vectors (comparable to Figure 8.23 in Ertel where per node in the tree the feature vectors are mentioned). For leaf nodes the DecisionTree contains a classification label (yes/no in Ertel's example, PLAY/DON'TPLAY in the Golf example of Figure 3). The DecisionTree class has several methods that still need to be implemented, for instance:

<sup>5</sup> Note that this label differs from the label in LabeledDataset. The latter one is a classification label, that is, it is the name of a class as it appears in the leaf of a decision tree.

```

2  // implementation of interface Classifier's method
    public L classify(F v) { ... }
4
6  /**
    * called from constructor of DecisionTree, this method should
    * do the actual recursive building of the decision tree.
    */
8  protected void build(LabeledDataset2<F, L> dataset) { ... }
10 public int prune(final LabeledDataset2<F, L> testData) { ... }
12 /**
    * compute the errorRate of this decision tree for testData.
    */
14 public double errorRate(Map<F, L> testData) { ... }
    
```