

AdaBoost

IT3105

Nordmoen, Jørgen H.
Østensen, Trond

November 19, 2012

Abstract

In this paper we will describe how we implemented AdaBoost, two weak classifiers and what results we got on a number of datasets. We will describe the datasets and why we chose them. We will explain how we divided those datasets into buckets for classification and how that is supported in our implementation. We will end this paper with a discussion of our results trying to explain and illustrate why.

Contents

1	Implementation	1
2	Results	2
2.1	Glass Dataset	2
2.2	Nursery Dataset	4
2.3	Pen Digit Dataset	6

List of Tables

1	Table showing the results of our classifiers on the Glass dataset .	3
2	Table showing the results of our classifiers on the Nursery dataset	5
3	Table showing the results of our classifiers on the Nursery dataset	7

1 Implementation

In this section we will talk a bit about our implementation. We will describe how it manages to classify new instances, how we generate new classifiers and how it can be extended without changing any implementation.

From our AdaBoost implementations view all we generate are Hypothesis¹ which comes from Generators². These interfaces dictates what all new classifiers must support. A Generator creates a new Hypothesis with the basis in weights and the training set. An Hypothesis only needs to support setting its weight and classifying a new instance³.

From an Hypothesis point of view all data is of type DataElement⁴ which contains a number of attributes specified when it's created. It also contains a classification which it must be able to produce. DataElement uses Java generics to support different types of attributes and classifications which is determined by the filter used. Each Generator must support the same attributes and classifications as the DataElement used, but this is only checked at runtime.

Our data flow starts with reading all lines in a file creating DataElements with String values. Then we apply a filter which must implement the interface Filter⁵. After that is done we shuffle the dataset using a specified random seed to be able to replicate our tests. We then split the dataset into a training set and test set with a split specified on the command-line. After that we start AdaBoost and try to classify the test set.

The meat of our implementation lies in the Generator class which dictates how a specific classifier is created. This is the class which is defined on the command-line and this is the class which can be configured from the command-line.

This class must be able to receive a list of options where it can act on each option according to its own wishes which is done so we can configure each Generator, this is among other things where we pass the depth parameter to the Decision Tree Generator.

Each Generator class uses its own logic to create a classifier of the correct type based on the given data. This means that AdaBoost does not know nor care what a Generator creates as long as it implements the Hypothesis interface. This also means that our implementation is very general. When creating a new classifier type one has to implement a new Generator and a new Hypothesis class, which is all. To support new datasets one only has to implement a new filter to convert from String to a specific type. This was done to decouple the datasets from the classifiers, but it also means that we have to implement less code since each classifier shares the filters.

The filters are tasked with both converting the String values read from a file to the correct format, but it's also tasked with creating buckets for the classifiers. With this we mean that the filter must split the values in each file into proper buckets since our classifiers can't handle large ranges. This is also the reason why we have a filter per dataset. Since we had a hard time creating

¹ no.ntnu.ai.hypothesis.Hypothesis

² no.ntnu.ai.hypothesis.Generator

³We have created an abstract class for convenience, situated in no.ntnu.ai.classifiers.Classifier

⁴ no.ntnu.ai.data.DataElement

⁵no.ntnu.ai.filter.Filter

an algorithm for splitting the ranges in different datasets we chose to have it like to for ease and to be able to control the buckets in more detail.

2 Results

Below is all the results from all our runs. We have used a couple standard values to control the randomness and the split of the datasets.

For each run we have used a split of "0.2" which means that we use 0.2 of our dataset as a test set. We have also used the random seed "42" for all runs this means that in every place where we need a random number the generator for that number is initialized in the same way with the same number each time. This creates reproducibility and ensures that we can be certain that our results are given the same starting point every time.

For each dataset we will point out what is needed to replicate the results.

2.1 Glass Dataset

This dataset represents the classification of different glasses in a crime scene.

To run, use these options with the classifier options wanted.

Listing 1: Glass dataset general options

```
—global 0.2 42 —file glass.txt —filter no.ntnu.ai.  
filter.GlassFilter
```

The classifier options below does not conform to our actual command-line style which is down to space requirement. To use just add the sentence ' — classifier no.ntnu.ai[nbc—dte].' in front.

NBC #	Training Error	Standard Deviation	DTC #	Training Error	Standard Deviation	Test Error	Classifier option
1	0.339	0.0	0	0	0	18/43(41%)	NBCGenerator 1
0	0	0	1	0.181	0.0	17/43(39%)	DTCGenerator 1
5	0.391	0.045	0	0	0	19/43(44%)	NBCGenerator 5
10	0.371	0.042	0	0	0	19/43(44%)	NBCGenerator 10
20	0.358	0.034	0	0	0	20/43(46%)	NBCGenerator 20
0	0	0	5	0.256	0.043	18/43(41%)	DTCGenerator 5
0	0	0	10	0.673	0.133	23/43(53%)	DTCGenerator 10 1
0	0	0	10	0.623	0.124	24/43(55%)	DTCGenerator 10 2
0	0	0	10	0.265	0.047	18/43(41%)	DTCGenerator 10
0	0	0	20	0.268	0.042	19/43(44%)	DTCGenerator 20
5	0.430	0.046	5	0.651	0.095	21/43(48%)	DTCGenerator 5 2, NBCGenerator 5
10	0.391	0.044	10	0.646	0.105	17/43(39%)	DTCGenerator 10 2, NBCGenerator 10
20	0.409	0.045	20	0.661	0.093	22/43(51%)	DTCGenerator 20 2, NBCGenerator 20

Table 1: Table showing the results of our classifiers on the Glass dataset

2.2 Nursery Dataset

This dataset represents the classification of applications for a nursery school.

To run, use these options with the classifier options wanted.

Listing 2: Nursery dataset general options

```
—global 0.2 42 —file nursery.txt —filter no.ntnu.ai.  
filter.NurseryFilter
```

NBC #	Training Error	Standard Deviation	DTC #	Training Error	Standard Deviation	Test Error	Classifier option
1	0.097	N/A	0	N/A	N/A	255/2592(9%)	NBCGenerator 1
0	N/A	N/A	1	0.0	N/A	76/2592(2%)	DTCGenerator 1
5	0.234	0.117	0	N/A	N/A	249/2592(9%)	NBCGenerator 5
10	0.193	0.094	0	N/A	N/A	249/2592(9%)	NBCGenerator 10
20	0.161	0.094	0	N/A	N/A	249/2592(9%)	NBCGenerator 20
0	N/A	N/A	5	0.0	0.0	76/2592(2%)	DTCGenerator 5
0	N/A	N/A	10	0.381	0.066	847/2592(32%)	DTCGenerator 10 1
0	N/A	N/A	10	0.358	0.081	451/2592(17%)	DTCGenerator 10 2
0	N/A	N/A	10	0.0	0.0	76/2592(2%)	DTCGenerator 10
0	N/A	N/A	20	0.0	0.0	76/2592(2%)	DTCGenerator 20
5	0.229	0.075	5	0.380	0.077	291/2592(11%)	DTCGenerator 5 2, NBCGenerator 5
10	0.194	0.075	10	0.349	0.079	313/2592(12%)	DTCGenerator 10 2, NBCGenerator 10
20	0.269	0.037	20	0.347	0.085	338/2592(13%)	DTCGenerator 20 2, NBCGenerator 20

Table 2: Table showing the results of our classifiers on the Nursery dataset

2.3 Pen Digit Dataset

This dataset represents the classification of digits written by different writers.

To run, use these options with the classifier options wanted.

Listing 3: Pen digit dataset general options

```
—global 0.2 42 —file pen-digits.txt —filter no.ntnu.ai  
    .filter.PenDigitsFilter
```

NBC #	Training Error	Standard Deviation	DTC #	Training Error	Standard Deviation	Test Error	Classifier option
1	0.119	N/A	0	N/A	N/A	251/2199(11%)	NBCGenerator 1
0	N/A	N/A	1	0.0	N/A	252/2119(2%)	DTCTGenerator 1
5	0.314	0.106	0	N/A	N/A	248/2199(11%)	NBCGenerator 5
10	0.329	0.081	0	N/A	N/A	239/2119(10%)	NBCGenerator 10
20	0.330	0.078	0	N/A	N/A	230/2119(10%)	NBCGenerator 20
0	N/A	N/A	5	0.0	0.0	252/2119(11%)	DTCTGenerator 5
0	N/A	N/A	10	0.761	0.057	1126/2119(55%)	DTCTGenerator 10 1
0	N/A	N/A	10	0.615	0.127	398/2119(18%)	DTCTGenerator 10 2
0	N/A	N/A	10	0.0	0.0	252/2119(11%)	DTCTGenerator 10
0	N/A	N/A	20	0.0	0.0	252/2119(11%)	DTCTGenerator 20
5	0.289	0.086	5	0.664	0.063	172/2119(7%)	DTCTGenerator 5 2, NBCGenerator 5
10	0.351	0.049	10	0.612	0.0121	196/2119(8%)	DTCTGenerator 10 2, NBCGenerator 10
20	0.347	0.073	20	0.656	0.091	174/2119(7%)	DTCTGenerator 20 2, NBCGenerator 20

Table 3: Table showing the results of our classifiers on the Nursery dataset