

Dirty Money:

Feature selection using AdaBoost

J. van Turnhout (0312649) jturnhou@science.uva.nl,
N. Raiman (0336696) nraiman@science.uva.nl,
S. L. Pintea (6109969) s.l.pintea@student.uva.nl

February 4, 2010

Abstract

In the month January of 2010 a project on the classification of the fitness of money was proposed. During this month we have tested and implemented various techniques to handle the problem of money classification. The results from our experiment show promising development comparing to the current state of research.

In the below sections we have described the approaches that we have tried and the results obtained for each one of them.

1 Introduction

The goal of this project is to research several techniques that have a potential to distinguish between clean bills and dirty bills. In order to achieve this goal we had to think of what are the representative features of dirt that can be found on money bills and what techniques can be applied to model this.

The main techniques used in this project are: *PCA*, *Haar-like features*, *Convolution with pre-defined kernels*, *edge-detection* and *intensity* distributions. All of these techniques have been analysed on the image data using a machine learning technique called Adaptive Boosting, also referred to in literature as AdaBoost [2].

In section 2, related work on the project is discussed. In section 3, the implementation of the *AdaBoost* algorithm applied on the different techniques *Haar-like features*, *convolution with kernels*, *PCA*, *SVM*, *AdaBoost*, *edge detection* and *intensity* are discussed. We discuss our experiments and their results in section 4. Finally, we conclude in section 5 and propose some topics

for future research.

2 Background

In Europe a lot of different factories produce euro bills using different ink and paper. This makes the task of detecting dirty bills hard because looking at the global features will not work perfectly due to these small differences in printing. In Holland, on average bills return 2 to 3 times a year to the Dutch National Bank (DNB) with a total value of 1.1 billion bills a year.

The current approach used in the DNB is to measure the reflection of the bill on a small area near the water mark. With this approach 30% of all bills are destroyed in order to destroy 95% of the dirty bills. In a first attempt to make this selection process more sophisticated global features (eigen-money) of clean and dirty bills are extracted using PCA analysis [3]. This approach improved on the results of the current method used by the DNB. The implementation was again improved by learning more local features of the water mark region (also referred to as white patch) of the bill [4]. This approach again improved the current results from DNB.

3 Methods

3.1 Adaboost

The *AdaBoost* algorithm will be used in combination with different techniques throughout this project, such as: *PCA*, *Haar-like features* or *edge* and *intensity* distributions over different regions.

AdaBoost is a machine learning technique which can be used in conjunction with various

other learning algorithms. The idea is to have a (convex) set of weak classifiers (classifiers that perform at least better than random) and then minimize the total error over the training-set by finding the best classifier at each stage of the algorithm.

The theoretical basis of this algorithm is that given a set of models (or features), M , the algorithm will determine the subset of T models that are the best for distinguishing between the two classes (clean and dirty bills). Thus, *AdaBoost* will learn the most representative features of the two classes. The algorithm for determining the best models is shown in *Algorithm 1*. Another very important characteristic of this algorithm is the fact that it also specifies a method in which the models that were chosen as being the best, can be combined in order to give a strong classifier. The corresponding algorithm can be seen in *Algorithm 2*.

Note that the algorithm looks slightly different but is equivalent to the one described in [1].

Algorithm 1 AdaBoost learning features

```

1: function AdaBoostLearn( $T, M, S$ )
2:  $T = \text{nr. of hypothesis}$ 
3:  $M = \text{Models}$ 
4:  $S = \text{training-set}, \{(x_1, y_1), \dots (x_n, y_n)\}$ 
   with  $x_i \in X$  and  $y_i \in \{-1, 1\}$ 
5:  $D_{1(i)} \leftarrow \frac{1}{n}$ , with  $i = 1, \dots, n$ 
6: for  $t = 1$  to  $T$  do
7:    $\text{error}_t \leftarrow 0$ 
8:   for  $m \in M$  do
9:      $h_j(x_i) \leftarrow \text{predict}(x_i) \% \text{svm or gaussian}$ 
        $\text{distribution}$ 
10:     $\text{error}_j \leftarrow \sum_{i=1}^n D_t(i)[y_i \neq h_j(x_i)]$ 
11:    if  $\text{error}_j < \text{error}_t$  then
12:       $\text{error}_t \leftarrow \text{error}_j$ 
13:       $h_t \leftarrow h_j$ 
14:       $\alpha_t \leftarrow 0.5 \cdot \log \frac{1 - \text{error}_t}{\text{error}_t}$ 
15:      for  $i = 1$  to  $n$  do
16:         $D_{t+1}(i) \leftarrow \frac{D_t(i) \exp(-\alpha_t \cdot y_i \cdot h_t(x_i))}{Z_t}$ 
17: return  $\alpha, h$ 

```

3.2 Haar-like features and convolution with Haar-like kernels

The *Viola & Jones* approach for object detection has been successfully applied to face recogni-

Algorithm 2 AdaBoost Prediction

```

1: function AdaBoostPredict( $\alpha, h, I$ )
2:  $\alpha = \text{weights}$ 
3:  $h = \text{weak classifiers}$ 
4:  $I = \text{image}$ 
5:  $p = \text{prediction}$ 
6: for  $t = 1$  to  $\text{length}(\alpha)$  do
7:    $p \leftarrow p + \alpha_t h_t(I)$ 
8: return  $\text{sign}(p)$ 

```

tion tasks. In previous projects, the use of PCA for classification is motivated by its application in face recognition. The use of Haar-like features has proven its purpose in this field as well and therefore its plausible the technique could be adapted to our problem. The final cascade used in the paper however was not implemented because the main idea of this cascade was not suitable for the purpose of this project. We have used the strong classifier computed in *AdaBoost* as the final output.

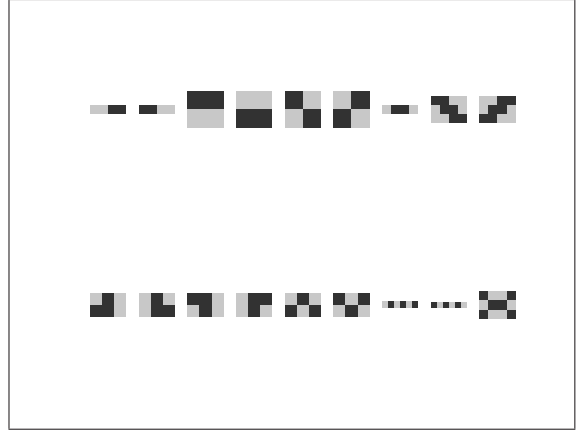


Figure 1: Patterns

The first step of the algorithm is defining the patterns, that are mainly matrices of different dimensions containing low and high values for intensity (-1 and 1). Figure 1 indicates a subset of the patterns used. The algorithm used in the *Viola & Jones* article was designed to loop over all predefined patterns on each position in the image and convolve the specific region of the image with the patterns using the formula depicted in figure 2.

We have started by defining a large number of patterns of different sizes (≈ 100 patterns), and filtered out those that gave bad results when

Figure 2:

$$\sum_x \sum_y Image(y : y + h, x : x + w) * Pattern$$

where: h – the height of the pattern
 w – the width of the pattern

used in the *AdaBoost algorithm* for training the weak classifiers. The resulted values for each pattern and location in the image would then be used in *AdaBoost* to train an *SVM* classifier.

Taking into account the fact that the set of features generated by the algorithm for each pattern and each image was extremely large, the training took too much time. In order to solve this problem, we have decided to use random locations at which to convolve the patterns with a region of the image that would have the same size as the pattern. Figure 3 indicates what *AdaBoost* would choose as being the most representative 7 features for the front side and rear side of the bills.

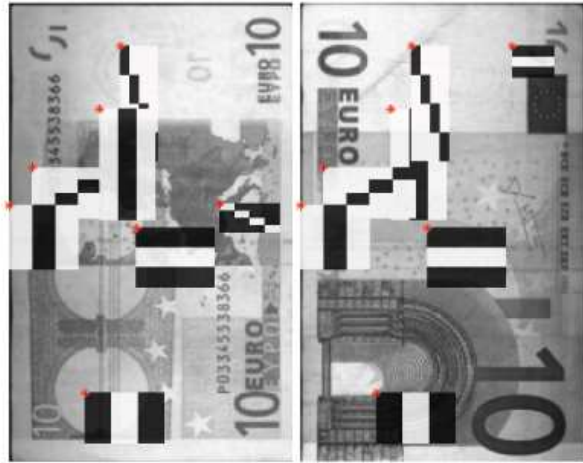


Figure 3: Haar Features for Front and Rear

The results obtained using *Haar features* were not as good as we would have expected them to be. An explanation for this may be the fact that the patterns defined were describing local features on pixel-level making it impossible to generalize to a more global level.

The second approach that we have tried was to define a set of patterns as before, and to segment each image into smaller regions that would be,

then, convolved with the patterns.

We have tried using both non-overlapping segmentation of the images and 50% overlapping one. Due to the fact that many essential pixel values (like those that would indicate the presence of dirt) may happen to be positioned on the line that separates two neighboring regions, it is possible that some information will be lost. Taking into account the number of regions generated for all images we noticed that there is a large amount of pixel values that might not be considered in the case of non-overlapping segmentation. As expected, the second method of dividing the images into regions gave better results compared to the first one.

The results obtained using this technique represented the input set of features for the *AdaBoost algorithm*. Figure 4 will help getting a better understanding of how the convolved regions of image and patterns would look like. In this image is plotted the result obtained when convolving a bill with a simple pattern such as: $[1 \ -1; -1 \ 1]$ with an entire bill.



Figure 4: Convolution of an entire bill

The resulted set would represent the input features used in *AdaBoost*. For establishing the set of the best T features we have tried using both *SVM* and *Gaussian distributions*, and the results retrieved by the last one seemed slightly better than the ones obtained while using *SVM*. In the case in which *SVM* was used, a model was generated for each available feature. This model would give the best separation between the values corresponding to that feature for fit images and those for unfit images. In the case of the *Gaussian distribution*, the *mean* and *covariance* of features corresponding to fit images were

computed, respectively the mean and covariance for the features corresponding to unfit images. In order to determine the predicted class of the input images we have tried two methods.

In the first method we would simply compute the *mean* and *covariance* for the values corresponding to fit, respectively unfit values corresponding to a specific feature. During testing, for each input value, we would compare the two numbers returned by *Gaussian density function* for fit, respectively unfit and we would choose the predicted class to be the maximum of the two.

The second method was using *Naive Bayes* in order to make use of the prior knowledge available. We know that in a real-life situation there are always more fit bills than unfit one and we would like to use this information to improve our classifiers. The predicted class was defined by using the *MAP* (maximum aposterior probability) estimation. The formula for computing the maximum aposterior probability for fit, respectively unfit class is given by the formulas depicted in figure 5.

Figure 5:

$$\begin{aligned} \operatorname{argmax}_{\theta_{fit}} P(\theta_{fit} | X) &= \frac{P(\theta_{fit}) * P(X | \theta_{fit})}{P(\theta_{fit}) * P(X | \theta_{fit}) + P(\theta_{unfit}) * P(X | \theta_{unfit})} \\ \operatorname{argmax}_{\theta_{unfit}} P(\theta_{unfit} | X) &= \frac{P(\theta_{unfit}) * P(X | \theta_{unfit})}{P(\theta_{fit}) * P(X | \theta_{fit}) + P(\theta_{unfit}) * P(X | \theta_{unfit})} \end{aligned}$$

where:

- $P(\theta_{fit}), P(\theta_{unfit})$ – marginal probabilities of "fit" class, respectively "unfit" class
- $\theta_{fit}, \theta_{unfit}$ – the parameters of the classes (mean and covariance)
- $P(X | \theta_{fit}), P(X | \theta_{unfit})$ – the conditional probabilities of the two classes (representing normal distributions)

It can be easily shown that the values of the parameters θ for a normal distribution (mean and covariance) that would maximize the probability of the classes are exactly the corresponding formulas indicated in figure 6.

The second technique used for determining the predicted class of the data given the parame-

Figure 6:

$$\begin{aligned} \mu &= \sum_i \frac{x_i}{|X|} \\ \sigma &= \frac{(X - \mu)(X - \mu)^T}{|X|} \end{aligned}$$

where X – the training data set
 μ – the mean
 σ – the covariance

ters seemed to provides slightly better results in practice due to the fact that it also incorporates some prior knowledge of the data.

For this approach we have used two classifiers: one for the rear side and another for the front side of the bills and the predictions given by the two classifiers were combined into a third one using Naive Bayes.

In order to create the final model, we have tried two different methods:

- the first one was essentially just choosing the best model (the one that gives the minimum error) throughout all the repetitions and all the rounds in the validation;
- the second method was to use a system of voting – each time a feature was chosen by the *AdaBoost* algorithm, it would receive a vote. In the end the top T most voted features throughout all the validation rounds and all the repetitions would be chosen from the best model found. The corresponding weights would be generate by taking the mean of the features' α -weights computed in the *AdaBoost* algorithm;

As we would have expected, the voting method for defining the best model, proved to give better results in this case so it was chosen to be applied.

3.3 PCA and SVM

In previous research, the use of SVM in combination with PCA already gave great error reduction on both the 5 and 10 euro bills. The best results were conducted from performing PCA on the whitepatch area only, and propagate the projections through the SVM. PCA on the whitepatch as well as the whole bill made it possible to reduce the data from a dimensionality equal to the number of pixels of the respective area to only

a 30 components. This reduction ensures faster training and classification with SVM. The idea is to identify more regions with similar discriminative features as the whitepatch. This is done by dividing the bill up in an arbitrary number of regions, and performing PCA for each of these regions in all the bills. The PCA segments then can be used to train SVM models for each segment. These SVM models in combination with their respective PCA segments in turn can be used as weak classifiers in AdaBoost. The dimensions for division are chosen using empirical evidence of the error rates, and as to be expected the segments cannot be too local for PCA to work.

3.4 Intensity and Edge distributions

Next to the more complex approaches we decided to implement also two more simple techniques, based on the intensities on the image and on the edges on the bill. The idea of using the intensity is based on the current approach used in the DNB (e.g. measuring the reflection of light on a small patch near the watermark). The idea of using the edges is based on the assumption that dirty and old bills have significantly more wrinkles and folds than relative new bills. In the first simple implementation the raw intensity values of the image of the bill are used. This is done by segmenting the bills into small regions (5 by 12 showed best results). Over these regions the average intensity is calculated. In this way it is possible to retrieve two Gaussian distributions of the intensity of all the regions over all the images, one for clean bills and one for dirty bills. This provides a tool to compare the probability of a region on a bill belonging to a clean or dirty bill.

In the edge approach the image is first filtered using the canny-edge detector. The canny edge filter works in 5 steps:

- Step 1: Gaussian smoothing
- Step 2: Extract the gradient of the image by convolving the image with 2 kernels, one to find the derivative in the X direction and one to find the derivative in the Y direction
- Step 3: Determine the direction of the edge using results from step 2

- Step 4: Apply nonmaximum suppression to find the real edge (selecting only the maximum edge points found)
- Step 5: Finally the hysteresis is used to make the lines continues

The result of an image of a bill after applying canny edge filter can be seen in figure 7.

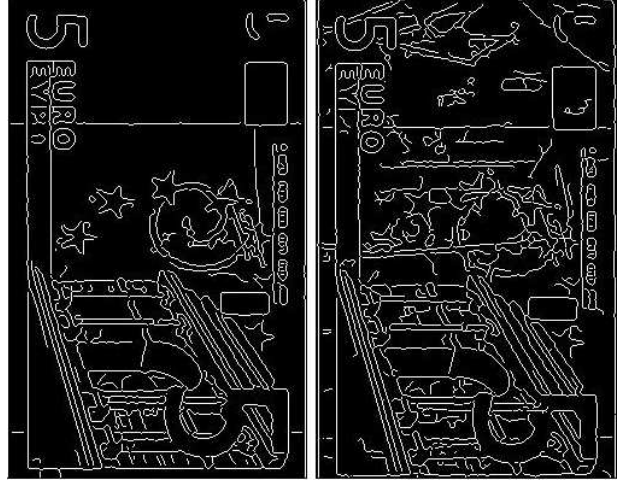


Figure 7: results canny edge filter. Left on a clean bill. Right on a dirty bill

After applying the canny edge filter to the images the resulting images are segmented in small regions. Per region the sum-of-edge-points is calculated. When doing this for all images we can again extract two Gaussian distributions of edge points per region (one for clean and one for dirty bills). Since calculating the average intensity and sum-of-edge-points only need to be calculated once, the learning phase is much faster than compared to the other two techniques which makes it possible to have more convergence in finding the best model during training.

4 Results

The available image data set was split into a holdout-set and the remaining data was used to develop and train our algorithms on. The holdout-set consists of 60 fit and 40 unfit bills, selected randomly out of the total set of images (both for 5 and 10 euro). During training on the remaining set, it was split up into a validation-set, a training-set and a test-set to perform *random sub-sampling validation*.

The whole process of defining the validation set and applying random sub-sampling validation was repeated several times to ensure a correct estimations of the error-rates. For each round of the validation, a model was trained using the *AdaBoost algorithm* described in section 3. The obtained model after training was tested by building the strong classifier and computing the corresponding confusion matrix values: *true-positive estimation* (unfit classified as unfit), *true-negative estimation* (fit classified as fit), *false-positive estimation* (fit classified as unfit) and *false-negative estimation* (unfit classified as fit).

In order to derive a number of hypothesis to use for AdaBoost, we have plotted the error-rates for each technique for a number of hypothesis. However, different decisions had to be made here for each technique since some take much longer to train.

4.1 Convolution Results

Figure 11 and figure 10 show how the error evolves depending on the number of models for the 10 Euro bills, respectively for the 5 Euro bills. Due to the fact that the number of features generated for only 21 patterns and a segmentation for each image of 9 (in the x-direction) by 23 (in the y-direction) overlapping regions gives a total of 4,347 features, the training part will be relatively slow for this method. In order to create reliable models we have used 5 repetitions, 10 hypothesis (features to be chosen by AdaBoost), 10 trials in the random sub-sampling validation method. The choice for 10 hypothesis is based on the fact that for 10 euro bills the error-rate got stable around this number, and using more hypothesis would result in time consuming learning trials.

In figure 9 are indicated the most voted regions for the 10 Euro bills in *AdaBoost* throughout all trials and all repetitions. We can notice that for the front of the bills the regions selected are mainly positions on the water mark area, while for the rear side of the bills, there are some regions considered as being the most informative on the white patch, but also the area in the middle is selected as containing discriminative information. This could be explained by the fold in the middle of the bills.

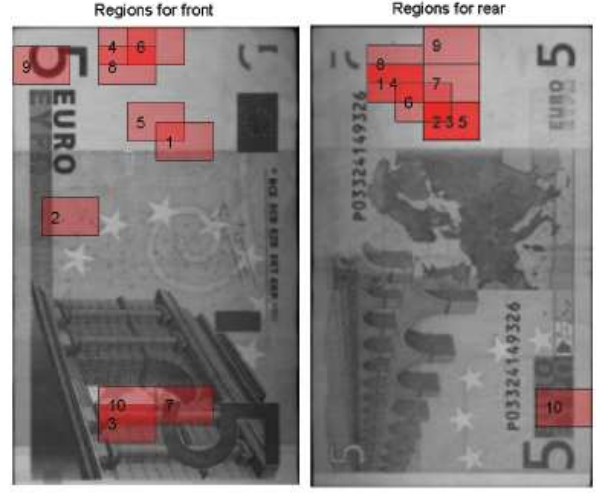


Figure 8: Best regions for 5 Euro bills (Haar)

Figure 8 shows the most voted regions for the 5 Euro bills in *AdaBoost* throughout all trials and all repetitions. For the 5 Euro bills the regions for front and rear that were chosen as being the most informative ones are again mainly around the water mark area, but unlike those for 10 Euro Bills, the regions for front side are the ones that are more spread throughout the image.

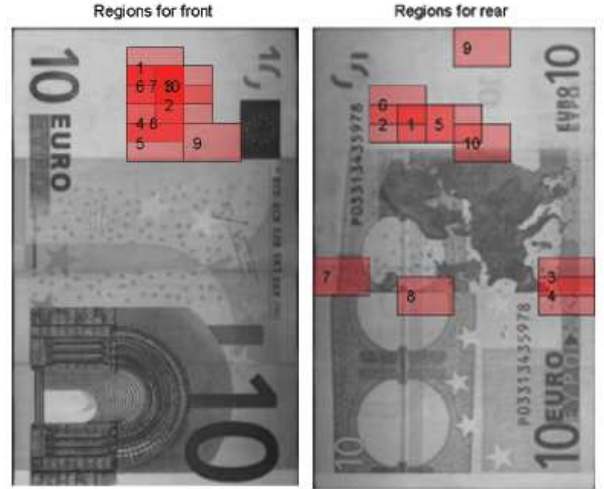


Figure 9: Best regions for 10 Euro bills (Haar)

4.2 PCA Results

Defining the number of regions to perform PCA on was derived from a few experiments on different number of segments. When trying too small segments, the information captured by

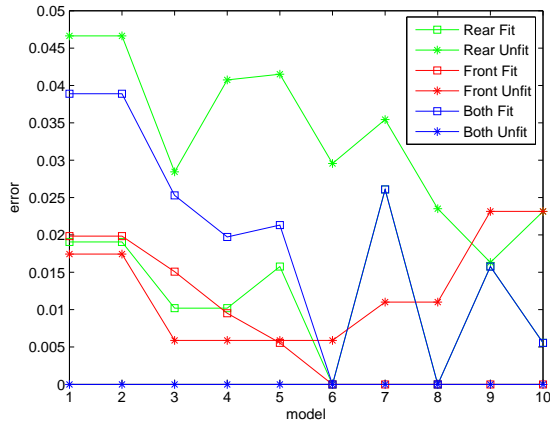


Figure 10: Error depending on the number of features - 5 Euro bills (Haar)

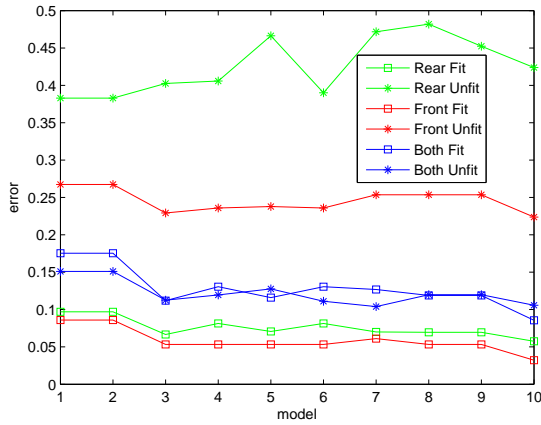


Figure 11: Error depending on the number of features - 10 Euro bills (Haar)

PCA could not represent the dirt on the bills anymore. When defining the segments too large, the information gets too general, and similar results are measured as for the whole bill. Dividing the bill into 3 (in the x-direction) by 7 (in the y-direction) regions gave promising error-rates during the training part as can be seen in figure 12 and 13. Surprisingly, the number of regions to choose to get a maximal error of 5% on unfit bills (the baseline for DNB) is only measured at 1 and 2 weak classifiers. Because the error for the fit bills is lowest when using one weak classifier only, the final model is trained on finding only one region. The most discriminating regions can be seen in figure 15 and ???. Note that the learning process of this technique is very slow because after obtaining a new train- and test-set during

random sub-sampling validation the eigen-faces for every region need to be calculated.

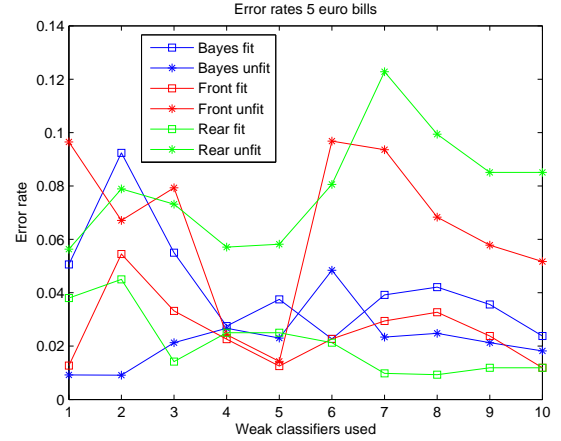


Figure 12: Error depending on the number of features - 5 Euro bills (PCA)

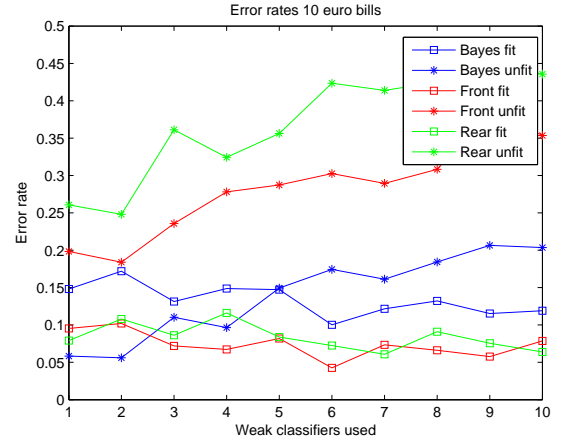


Figure 13: Error depending on the number of features - 10 Euro bills (PCA)

4.3 Intensity & Edges Results

For the intensity and edges implementation we also used overlapping regions, 9 (in the x-direction) by 23 (in the y-direction) for front and for the rear of the bill. This results in 414 weak classifiers/features. For every region there is only one value/number that represents the edge count or the average intensity, So there is a relative small amount of models to train on. Therefore we could run bigger tests in order to achieve reliable models. In order to create these models we have used 40 repetitions, 40 hypothe-

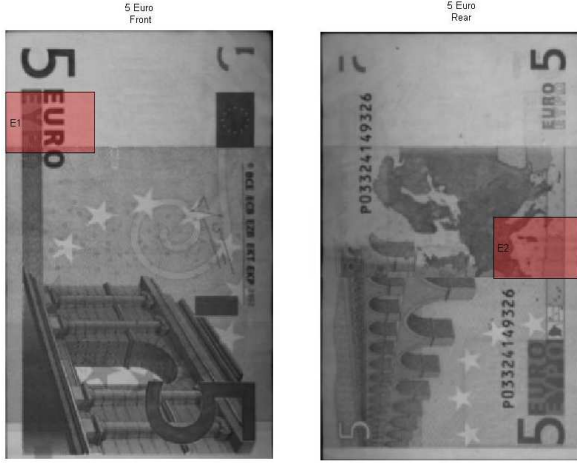


Figure 14: Best regions for 5 Euro bills (PCA)

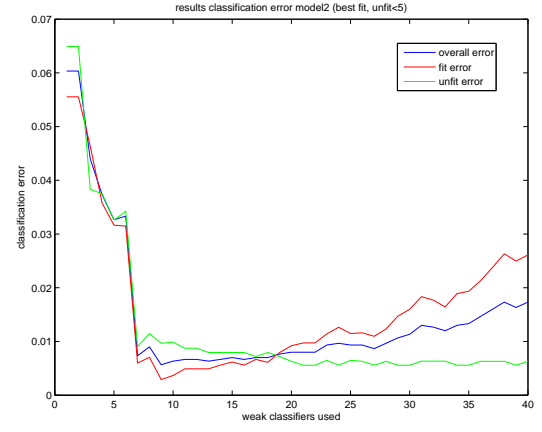


Figure 16: Error depending on the number of features - 5 Euro bills (IE)

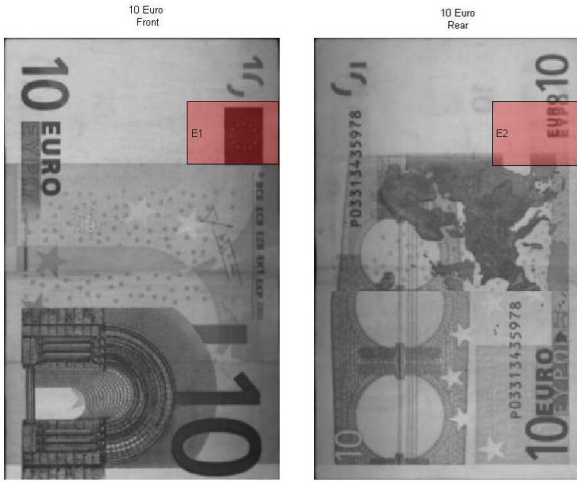


Figure 15: Best regions for 10 Euro bills (PCA)

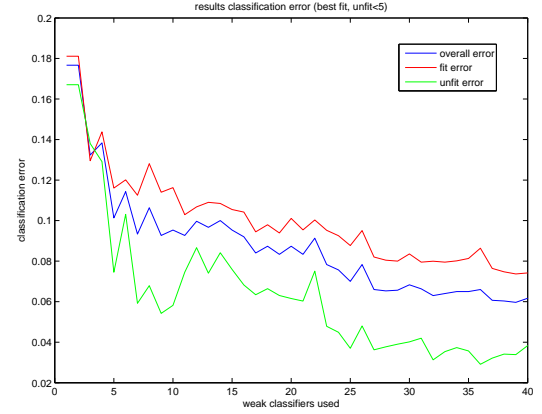


Figure 17: Error depending on the number of features - 10 Euro bills (IE)

ses and 40 random sub-sampling validation runs. To determine the optimal number of weak classifiers to use for the 5 and 10 euro bills we have plotted the errors with respect to the number of weak classifiers used. We choose the weak classifiers that reduced the error 'substantially' (see figure 16 and 17 10 and 5 euro plots). For the 10 euro bills we used 25 weak classifiers and for the 5 euro we used 10 classifiers (see figure 18 and 19 figures of areas on bills for 10 and 5 euro).

4.4 Combined Results

Once all the strong classifiers for each of the methods above have been generated, we have combined their predictions using *Naive Bayes* rule. Since all of our final models are prepared

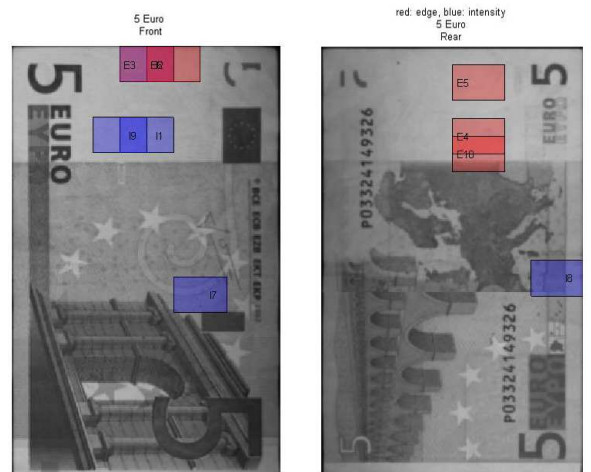


Figure 18: Best regions for 5 Euro bills (IE)

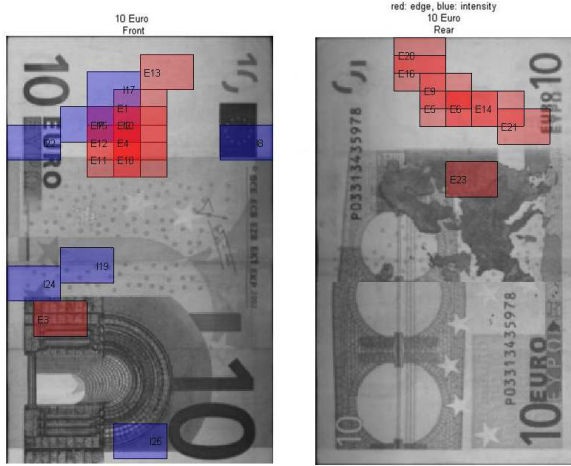


Figure 19: Best regions for 10 Euro bills (IE)

with a number of hypothesis which are equal or below the 5% baseline from DNB, we choose to bias the Naive Bayes towards the fit bills with the motivation that there are more fit bills being destroyed relative to the number of unfit bills. The combined classifiers are tested on our holdout-set of 100 bills for each 5 and 10 Euro.

The results for 10 Euro bills, respectively 5 Euro bills are shown in the Table2 and Table1. For the 5 Euro bills the performance of both single classifiers Haar and IE look very promising already, and in combination we can only expect them to be as good or slightly worse in performance.

Table 1: Results for 5 Euro bills

	<i>Fit Error</i>	<i>Unfit Error</i>
<i>Haar</i>	0	0
<i>IE</i>	0.033	0
<i>PCA</i>	0.083	0.025
<i>Haar & IE</i>	0	0
<i>Haar & PCA</i>	0	0.025
<i>IE & PCA</i>	0.033	0.025
<i>Haar & IE & PCA</i>	0.033	0

5 Conclusion

From our experiments we have been able to learn the regions that contain the largest amount of information and help us distinguish between fit and unfit bills. These regions differ from one method to another, but in general the areas

Table 2: Results for 10 Euro bills

	<i>Fit Error</i>	<i>Unfit Error</i>
<i>Haar</i>	0.05	0.15
<i>IE</i>	0.1	0.025
<i>PCA</i>	0.083	0.05
<i>Haar & IE</i>	0.033	0.15
<i>Haar & PCA</i>	0.017	0.2
<i>IE & PCA</i>	0.017	0.075
<i>Haar & IE & PCA</i>	0.033	0.025

around the water-mark region and the middle of the bills have been proven to be the ones containing the most discriminative features.

As it can be noticed from the results in section 4.4, combining different techniques can lead to an improvement in the performance of the final classifier although running single classifiers only give great results as well. In terms of robustness it make sense to do combinations of strong classifiers for 5 and 10 Euro both bills, since the different techniques can cover different aspects of dirt on the bills. To test this hypothesis however, it is mandatory to have many more holdout-sets to test the final models on since our project does not give any proof on scalability yet.

Although the results obtained using the various techniques discussed are very promising already, future work is possible and it should be mainly focused on finding a more powerful way of combining all the features used by all 3 techniques (Haar, PCA, Intensity & Edge) into a strong classifier as AdaBoost can be capable of doing this. Also, combining the models for both front and rear can also be improved if certain prior knowledge can be incorporated.

We would also recommend that special attention should be paid to the validity of certain features of the images from the data set by retrieving the origin on how the data-set is built up. We mention this because the results for average intensity are really promising, but this could be caused by bills being scanned in each class at the time and having the scanner being heated up.

References

- [1] P. Viola & M. Jones:
Rapid Object Detection using a Boosted Cascade of Simple Features (CVPR 2001)
- [2] AdaBoost:
<http://en.wikipedia.org/wiki/AdaBoost>
- [3] G. Molenaar, A. Nusselder and K. M. Stefanov:
Classifying Bank Notes for Scheduled Destruction (2009)
- [4] J.M. Geusebroek:
Computer Vision for Banknote Sorting (2009)