

Improving Performance and Applying Cascades Visual Classification

Master Thesis

Michal Kiwkwowitz

Contents

1	Introduction.....	
2	Evaluating Classification.....	
2.1	Evaluation Methods.....	
2.2	Practical Application of Classification.....	
2.3	Machine Classification vs. Human Classification.....	
2.4	False Negative Errors.....	
2.5	False Positive Errors.....	
3	Existing Classification Methods.....	
3.1	Feature Representation.....	
3.1.1	Fragments.....	
3.1.2	Gabor features and wavelets.....	
3.1.3	SIFT.....	
3.1.4	Improving Features - Hierarchies and Semantics.....	
3.1.5	Constellation Models.....	
3.1.6	Satellites for Similar Classes.....	
3.2	Learning and Classification.....	
3.2.1	Max-Min Selection and Bayesian Classification.....	
3.2.2	SVM.....	
3.2.3	Ensemble Methods.....	
3.2.4	Bagging.....	
3.2.5	Boosting.....	
3.2.6	Cascades.....	
3.3	Online Classification.....	
4	Improving Classification.....	
4.1	Learning More from the Training Set.....	
4.1.1	Improving the Similarity Measure.....	
4.1.2	Improving the Combination of Features into a Score.....	
4.1.3	Using the Configuration to Improve Classification.....	
4.1.4	Perspective Fragments.....	
4.1.5	Anti-Fragments.....	
4.1.5.1	Motivation and Algorithm.....	
4.1.5.2	Experiment and Results.....	
4.1.6	Satellite Fragments as a Second Stage.....	
4.1.6.1	Algorithm and Motivation.....	
4.1.6.2	Experiment and Results.....	
4.2	Classifying in the Setting of a Growing Training Set.....	
4.2.1	The Training set Limit.....	
4.2.1.1	Training set Experiments.....	
5	Cascades in Classification.....	
5.1	Three-Way Cascades.....	
5.1.1	Cascade Processing Region and Cost.....	
5.1.2	Orthogonal Cascade Processing.....	
5.1.3	On-Line cascades.....	

5.1.4	Three-Way Cascade Results	
5.2	Configuration Cascade	
5.2.1	Theory	
5.2.2	Algorithm	
5.2.3	Results	
6	Conclusions	
7	Future Work	
7.1	Improving Classification	
7.2	Improving the Anti Fragments concept	
7.3	Extending the Cascade	
7.4	Human Vision Experiments	
8	Appendix	
8.1	Max-Min Fragment Selection and Bayesian Classification	
8.2	AdaBoost	Error! Bookmark not def
8.3	Anti-Fragments	
8.4	Sattelite Fragments as a Second Stage	
8.5	Three-Way Cascade, Determining the Processing Region	
8.6	Online Three-Way Cascade using AdaBoost	
8.7	Configuration Cascade	
9	Acknowledgment	
10	References	

ABSTRACT

Computer-based object classification has improved consistently over the last decade, the performance of current computational schemes is still significantly lower than human classification performance. In addition, the errors made by current classifiers are unreasonable by human standards. Since improvements to performance become increasingly difficult to achieve when absolute performance is already high, it is unclear which future directions will be useful for reaching truly human level performance. In this work I examined two general directions for future improvements in current classification scheme. One general direction assumes that current methods extract only a part of the available information for classification from the training set, and attempt to identify the main possible sources of additional information. I examined a number of plausible sources for possible improvements. Some of the methods under study, but concluded that they are unlikely to be sufficient by themselves to reach human-level performance.

The second general direction assumes that a major limitation comes from the size of the training set: training sets used in practice may be inherently insufficient to capture all necessary variations needed to learn a truly high-performance classifier. Our experiments suggested that classification performance indeed increases monotonically and without clear saturation as the training set increases in size. We also found that the set of features used for classification needs to be increased to capture the additional information in the increased training set.

These results have two implications to the classification scheme. First, they raise the advantage of constructing classifiers in an on-line manner, in which the classifier is continuously improved by new examples. Second, they raise the problem of increasing computational load as the number of features used for classification increases.

To deal with these problems we developed and tested two classification schemes. Both multi-stage, or 'cascade' methods in which all images are analyzed in the first stage, depending on the results, only some of them are analyzed further by subsequent stages. The first of these methods was a so-called 3-way cascade, which is an extension of previously used cascade methods. In this method, a first-level classifier is first applied to the new image to be classified. If the response is high or low enough a decision is made, otherwise, a second-level classifier, which uses more features, is applied to the image to allow a more confident response. The second multi-stage scheme we developed is the so-called configuration cascade. In this approach, the decision to process an image through additional stages is based on the particular feature configuration discovered in the image. Using this scheme, erroneous configurations of the first-level classifier are processed by the second-level classifier which uses more features and achieves better performance. We showed that the number of possible erroneous configurations is bounded in practice, that the error on the erroneous configurations is reduced after additional processing. This scheme leads to a continuous improvement in performance, with only a small increase in computational cost.

1 Introduction

The field of visual recognition has been making continuous progress in the area of classification and recognition. However, up to date, no classifier for a real complex task has achieved classification without errors or fully human-level performance. In this work we examine the sources of these errors, and possible methods to improve performance without a large increase in computational cost.

We start by investigating the type of errors made, using as an example the fragment based approach [2], to classify faces and non faces, and evaluate in Section 2 the classification performance and the errors it makes. We then survey existing classifiers and their advantages and drawbacks, in Section 3. In Section 4, we examine different ideas for improving the performance of the fragment based classifier. We continue by identifying the training set size as an inherent boundary to the performance, and the implications of this observation on classifier's computational cost in Section 4.2. We suggest in Section 5 two methods of cascades, which can improve performance with relatively low additional computational cost. The first, a three-way cascade described in Section 5.1, performs further processing on images near the margin, and the second, a configuration cascade described in Section 5.2 which uses the configuration of the detected features to determine whether an image requires additional processing in the cascade.

In the next section we examine more closely the performance of current classifiers and their goals and motivations for the current work.

2 Evaluating Classification

The goals of our work are to answer the following questions:

- Can the performance of current state-of-the-art visual classifiers be considered satisfactory?
- Are we learning all we can learn from the training data?
- What are possible ways of improving performance further?

In this section we survey methods of evaluating classification performance. We consider the performance for practical applications and examine the amount of error

made by machine classifiers compared with human performance. The errors can generally be divided into two types: false negative errors, images that have been rejected falsely, and false positive errors, images in which the class has been falsely detected. We begin with an overview of the methods used to evaluate a classifier.

2.1 Evaluation Methods

Judging the performance of a classifier is not a trivial question. Many authors plot and compare the Receiver Operating Characteristic (ROC) curve (seen in Figure 1) of different classifiers, and although error rates may be small and satisfactory, the error instances themselves are quite alarming. Figure 2 and Figure 3 show some examples. This gives an understanding of the relationship between the different errors made by a classifier. A threshold on the response of the classifier determines a desired hit rate and respectively a false positive rate that depends on the classifier's ROC curve. In an ideal world we would be able to set the threshold so that our classifier will get 100% hits and 0% false positives, a perfect classifier with no errors. In reality, there is no perfect threshold, and the possibility is to select some acceptable mixture of false positives and false negatives.

The ROC, although very useful, is not a simple parameter to compare. A measure that is easier to compare is the Equal Error Rate (EER) of a classifier. This is the location on the ROC curve where the percent of false positives equals the percent of false negatives. For a perfect classifier the value of the EER will be zero and therefore there will be a threshold that allows the classifier to make a perfect decision every time. A similar measure is the Minimum Error which provides the minimum total error (false positives and false negatives). The EER and the Minimum Error are usually found at close locations on the ROC curve. Another way to appraise the ROC is to measure the percent of the area below the ROC curve; the perfect ROC will have a value of one.

For most real world classification problems, there are no perfect classifiers; no classifier improves the classification by improving the ROC curve at the top left corner. In many cases the changes are very subtle and the comparison of EER, area or even different ROC curves misses a lot of information about the classifier. Given two classifiers with the same ROC or very close ROC how can they be compared? Consider a vendor

is purchasing a classifier that uses images to alert him when a police car passes by property. He is offered two classifiers with very good ROC curves but imperfect c Shouldn't he consider the type of false positives made by these as an additional method rating their performance? Wouldn't a classifier that detects a cab as a police car e once in a while be better than one that detects a cat or just the wind in the trees? Simil when a child calls a cow a horse, his parents will happily just let him know that I wrong, but shouldn't they pay a little more attention if their child is using the word when he encounters shoes or when staring at a clear wall? The error instances themse are important and there is much to learn from them. We next examine the EER and type of errors made by current classifiers.

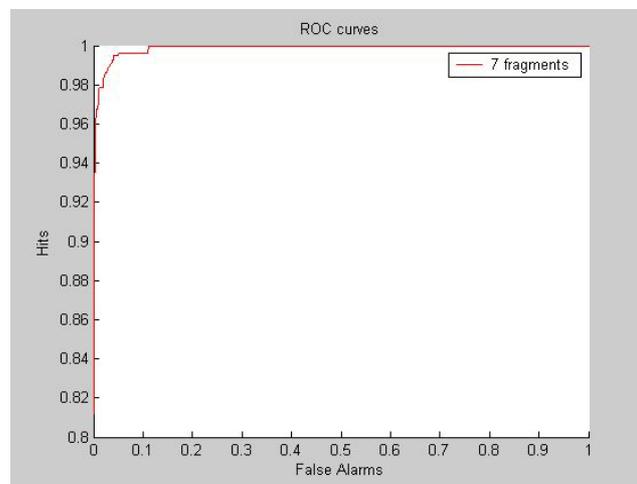


Figure 1: Classifier Performance. The Receiver Operating Characteristic curve of a classifier using seven fragments as described in [3]. Although the ROC is good, imperfect.

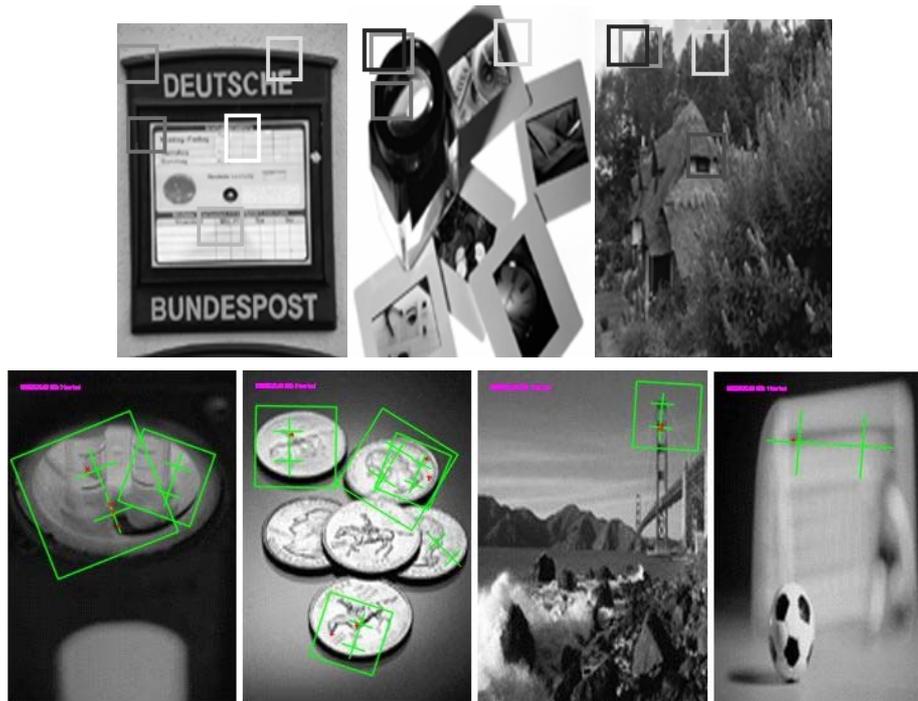


Figure 2: Faces Found by Face Classifiers. Images detected as faces by two different face classifiers with high performance. On the top row are results from a fragment classifier as in [3], the fragments detected are marked. On the bottom row are false faces detected by the Betaface classifier [8]. This classifier detects the eyes in an image and marks them in green.

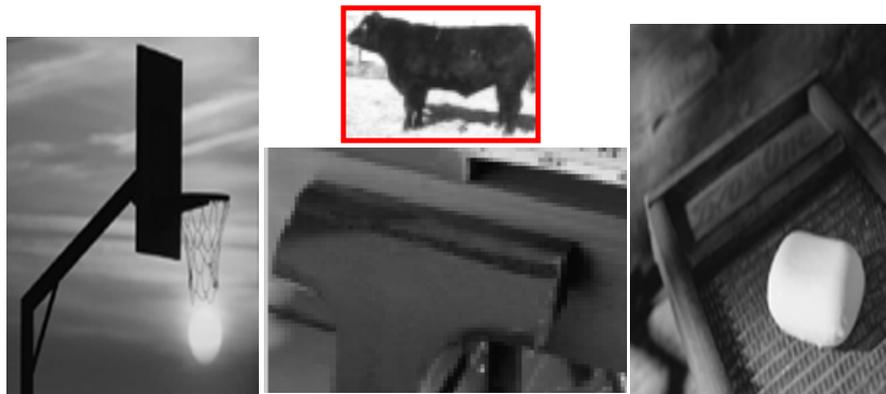


Figure 3: Cows Found using SIFT. These images received very high scores and were classified as cows by a classifier that uses SIFT features [7], we show a sample cow in a red rectangle.



Figure 4: **Classifier Performance on Specific Images.** On the top row are some missed by the classifier, in the center are some false alarms or false positives. The bottom row shows some images that we believe may serve as more reasonable false alarms.

2.2 Practical Application of Classification

Although the equal error rate in some of the best classifiers can reach 2-6 %, this is still quite high. Consider searching a natural scene and trying to determine if there is a horse in the image. Suppose that the image is 1000 pixels high by 1000 pixels wide, and we expect the horse object to be found in a window of around 50 by 50 pixels. Searching the image serially with local windows with a 50% overlap, we will need to check 40x40 locations, which is 1600 sub windows. With an error rate of 2%, this means that 32 sub windows will be falsely detected as horses, which is not very practical. In the next section we contrast human performance with the performance of machine classifiers.

2.3 Machine Classification vs. Human Classification

If humans had the false alarm rate mentioned in the previous section, they would hear horses everywhere, or in other words they would be prone to mistakes and visual hallucinations. Compared with State-of-the-art classifier's performance, humans perform better: for most databases or classification tasks there is no question whether machine performance is better or poorer than that of humans, humans outperform current classifiers. When comparing machine performance to human performance, the task classification is harder on machines, given that the in-class variability is greater. Machine software is closer to human performance in somewhat simpler tasks like recognition, where conditions are controlled for rotation and lighting. Not many real comparisons have been made between human and machine categorization, perhaps because it is common knowledge that humans are better. A comparison study by Adler and Schuckers (2006) compared automatic face recognition technologies available in 1999, 2001, 2003, 2005 and 2007 with human performance and showed that as automatic face detection improves, average human performance is no longer better than the machine performance. This is impressive, yet we need to consider that some human performance is nevertheless better than the best automatic face recognition, and the task at hand is answering 'same' or 'different' to pair of images from the NIST Mugshots Identification Database, where a match is between images of the same person at different age points which may cause facial differences, see Figure 5 and [1] for more details. The field of face recognition classification is one of the better studied classes, on other classes (planes, motorbikes, dogs, etc.) performance is usually poorer. As mentioned in the previous section, humans do not make as many false positive errors as machine classifiers; this fact is important for understanding the human visual system and for improving machine classification. In the next section we take a closer view at some classification errors dividing them into false negative and false positive errors.

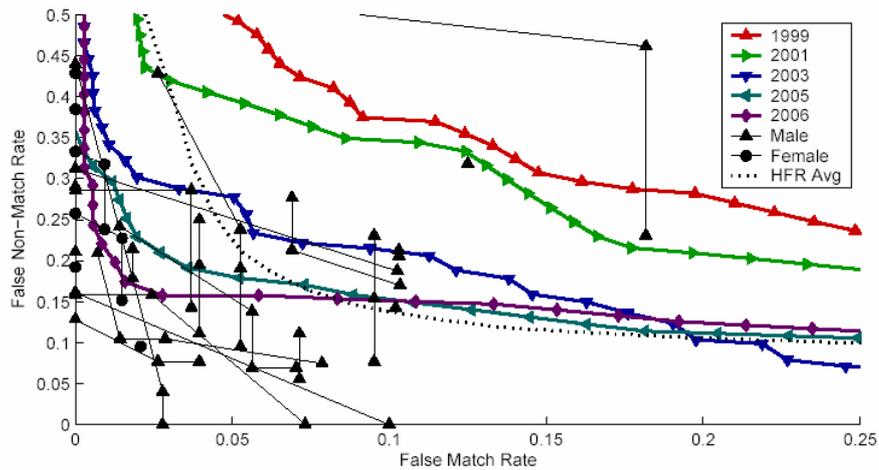


Figure 5: Comparison of face recognition software and human performance from The black circles and triangles indicate the human female and male perform respectively, and because they could respond on a scale of 1-5 conveying their confidence in the match, a curve was drawn for each human by adjusting the threshold on confidence, the continues curves show the results for the highest performing soft available to the authors of [1] in the specified years, and the dotted line is an average the human performance, that is highly affected by the outliers. Most human performance remained at the bottom left corner with the lowest error rates, and since where the human curves do not completely outperform an automatic curve, the authors counted it as an intermediate performance (33.3% comparing to the 2006 condition), they got a good score for the 2006 comparison. This may be challenged if the statistics were different.

2.4 False Negative Errors

When we examine the false negatives (misses) that current face classifiers make we often understand them (see Figure 4): many times misses occur in challenging face images with long hair, high foreheads, glasses, or in images that contain non frontal faces, shadows or noise. This seems intuitive and is consistent with human performance: infants are known to sometimes burst in tears when seeing a familiar kin that has grown a beard or put on new glasses. We can assume that large changes in a face that were not seen before may not be encoded as part of the stored features and may cause unfamiliarity. True errors are usually understandable; furthermore, algorithms have been suggested to reduce the rate of false negatives by adding more training samples, selecting more features, or boosting the classifier [13]. The false positive errors, discussed in the next section, are more interesting.

2.5 False Positive Errors

False positives errors are those images that the classifier falsely detects as class but are not. They can be divided into two main types of errors, which I call multi-class errors and single-class errors. As multi-class errors we consider those errors that are caused by confusion between similar classes. For example, when trying to classify horses, we sometimes recognize some cows as horses. Since both classes are four-legged animals, we consider these errors reasonable; extending the classifier into a multi-class classifier to learn similar classes, or adding the similar class to the training samples labeled as non-class can resolve this problem. Similar errors are also made by human infants before they learn the name of the new class.

The single-class errors are those images that do not resemble the class and are nevertheless classified as the class; see Figure 4 for some examples. Although the classifier shows a fairly good ROC curve, it makes unreasonable mistakes that a two-year-old would not make, failing to see for instance that a register is not a face. This problem of unreasonable false alarm (compared with human judgment) is not easily avoided, it persists in many different types of classifiers and when using different kinds of features. Using a Naïve Bayes classifier that extracts SIFT features as in [7] the problem persists, examples of their likelihood scores can be seen in Figure 6. In Figure 2 we show another algorithm, Betaface [8] that detects faces by finding the eye region, and fails as well when provided with queries of some non-class images that are hard for the fragments classifier.

The single-class false positives are interesting since it is clear they are unreasonable.

We continue to discuss some popular classification methods, their relative advantages and shortcomings in section 3. In Section 4 we use our observations to suggest and compare possible general methods for potential improvements in performance of the classifiers.



Figure 6: Cow Database Likelihood. On the top some examples of the likelihood scores produced by a cow classifier based on SIFT features [7]. The false positives are irrelevant to the class of cows. On the bottom some examples of images that we consider more relevant as false positives.

3 Existing Classification Methods

In this section I survey the main relevant classification methods, the features they use, how they combine features to reach a decision. The goal is to identify possible reasons for failure and possible methods for improvement.

Classification is the process of determining a label Y_i for each sample X_i for some unknown function $F: X_i \rightarrow Y_i$.