

# Training of multiple classifier systems utilizing partially labelled sequences

Martin Schels, Patrick Schillinger, and Friedhelm Schwenker

Ulm University - Department of Neural Information Processing  
89069 Ulm - Germany

**Abstract.** Making use of unlabeled data samples in a classification or training of a classifier is a desirable aim for many real world applications in pattern recognition. In this study, a multiple classifier system is utilized to investigate this matter. Further, cluster analysis is used in order to group the available data while neglecting the actual labels. Then, by implementing an information fusion architecture based on these clusters, a classification architecture is constructed. This kind of an architecture is investigated by means of a facial expression data collection with focusing on one-against-one class decisions to produce locally “unlabeled”, i.e. not assigned to one of the considered classes, data.

## 1 Introduction

Combining the classification powers of several classifiers is regarded as a general problem in various pattern recognition applications [1]. Several experimental and analytical investigations on *static* and *trainable* fusion schemes have been made in the literature (see for instance [2]). We consider the case where the individual classifiers are trained on predefined feature subspaces, such that the proposed multiple classifier system (MCS) shows a one-to-one correspondence between features and classifiers. A trainable fusion mapping is constructed by an extra supervised learning phase. Thus, training a MCS can be considered as a two phase learning approach with an architecture having two layers.

A major goal of this study is to analyse the effect of incorporating unlabelled data in the training procedure of individual classifiers of the first level. This is achieved by grouping the training set into clusters without taking into account the actual class labels. However, by constructing an information fusion layer, the team of individual classifiers can still be combined to yield a particular class label. Thus, one could be able to gain further classification stability by exploiting the unlabelled data than by just considering the labelled ones. Furthermore, in many real world applications, the amount of labelled data at hand is small, making this issue an interesting matter of research. Also, it might be convenient for a particular problem to neglect parts of the given label system of a particular application's domain and consider only a particular subspace. Nonetheless it

---

This paper is based on work done within the “Information Fusion” subproject of the Transregional Collaborative Research Center SFB/TRR 62 “Companion-Technology for Cognitive Technical Systems” funded by the German Research Foundation (DFG). The work of Martin Schels is supported by a scholarship of the Carl-Zeiss Foundation, and DFG project under contract SCHW623/4-3.

would be desirable to make use of all the available data samples. In addition, we hope to gain flexibility from this setting according to a change of the underlying label system, e.g. by relabelling the data by experts. In this paper such a situation is studied on the basis of an application to the recognition to facial expressions.

## 2 Supervised learning of sequential data

The most prominent technique for the processing of sequential data is the hidden Markov model (HMM) [3]. In this approach, the sample sequences are treated as observations, that are emitted from a latent stochastic process. Hence a HMM is defined as  $\lambda = (A, B, \pi)$ , where  $A = (a)_{ij}$  reflects the probability to pass from state  $i$  to state  $j$ ,  $\pi$  denotes the initial state probabilities and the probabilities to emit a particular observation for each state is subsumed in  $B$ . In case of a continuous multivariate observation space, the probabilities  $B$  are represented by a Gaussian mixture model (GMM) per state.

This construct poses 3 main problems: firstly, suppose a given sequence  $O$  and a HMM  $\lambda$ , now, what is the a posteriori probability of  $\lambda$  having emitted  $O$ ? This issue is dealt with the forward algorithm, introducing the forward variables, that are passing messages forward in time. Secondly, what is the most probable sequence of states in  $\lambda$  that has to be passed in order to emit  $O$ ? This matter is solved using the Viterbi algorithm. The hardest problem in this context is the estimation of adequate parameters of  $\lambda$ , given a set of sample observations. This training of a model is achieved by applying an expectation maximization algorithm: the Baum-Welch algorithm. This algorithm makes use of forward and backward variables, which are analogue for message passing backwards in time.

In order to conduct classification utilizing HMM, one model is created for each of the particular categories by making use of the respective available training data. For any new sample, that needs to be classified, every model is evaluated and the class label of the model that shows the highest a posteriori probability is assigned. However, it may be convenient to preserve the probabilities of all categories, in particular when these results are further processed later on.

## 3 Unsupervised learning considering sequences

In many real world applications, the class labels of distinct categories or even the categories itself may not be defined. Therefore, it is compelling to learn clusters of a data collection based of the similarities of the samples.

Because in general, considered sequences show different lengths (in contrast to stationary samples), a distance measure like the euclidean distance is not applicable in this context. Hence, a distance measure based on Hidden-Markov-Models as described [4] in is used. A distance-matrix for a set of sequences  $M = \{S_1, \dots, S_n\}$  is constructed using the following algorithm [5]:

1. Train one HMM  $\lambda_i$  per sequence  $S_i \in M$

2. Calculate the log likelihoods  $L_{i,j} = \log P(S_i|\lambda_j)$  for every sequence wrt. every HMM. In order to mitigate the effects that caused by the duration of a sequence, the log likelihood is normed using its length.
3. Finally the distance between two sequences is  $d(j,i) = \frac{1}{2}(\bar{L}_{i,j} + \bar{L}_{j,i})$

The result of this algorithm is a  $n \times n$  symmetric distance matrix, which is used as input to a hierarchical cluster analysis procedure. In general, any kind of hierarchical clustering approach is feasible. In this paper Ward's algorithm, minimizing the squared euclidean distance measure in each cluster fusion step [6], has been selected.

The idea of this work is to train classifiers with respect to the new label system, which is discovered by clustering. These classifiers must provide a confidence for the class memberships, e.g. the log-likelihoods provided by HMM. Based on these confidences, a supervised classification step can be conducted, e.g. as described in Sect. 4.

#### 4 Multiple classifier systems and classifier fusion

Instead of finding the best individual classifier for a particular task, another frequent strategy to solve a classification problem is the philosophy of multiple classifier systems (MCS) [2]. This implies to create more than one individual classifier and then, combine them in an appropriate way. There is extensive research on how and under what circumstance it is sound to apply classifier fusion [2, 7, 8, 9, 10]. Major findings imply that the individual classifiers need to be individually accurate but also diverse with respect to the classifier team [7].

There are different approaches to perform classifier fusion: in many applications, a fixed combination rule is applied [8], e.g. averaging classifier outputs or performing majority voting, however, it is also convenient to construct a further fusion layer from the available data (i.e. a trainable mapping from the individual classifier's output to the true label of a sample) [9, 10].

In this study decision fusion is conducted using a linear pseudo inverse mapping as described in [10]. Hence, a mapping matrix  $W$  is calculated from the output of the models to match the true label of a sample. Let  $Y \in R^{m \times p}$  a matrix comprising the  $p$  available vectorized labels of dimension  $m$  of the training set. Further,  $C \in R^{m \times n}$  denotes the output generated by the individual models. The pseudo inverse mapping [11] is then defined as:

$$X^+ = Y \lim_{\alpha \rightarrow 0_+} C^T (CC^T + \alpha I)^{-1}. \quad (1)$$

Here,  $I$  denotes the identity matrix and the superscripted  $T$  transposes the respective argument.

#### 5 Data Collection

The Cohn-Kanade dataset is a collection of image sequences with emotional content, which is available for research purposes [12]. It contains 432 image

sequences, which were recorded in a resolution of  $640 \times 480$  (sometimes 490) pixels with a temporal resolution of 30 frames per second. Every sequence is played by an amateur actor recorded from a frontal view. The sequences always start with a neutral facial expression and end with the full blown emotion which is one of the six categories “fear”, “happiness”, “sadness”, “disgust”, “surprise” or “anger”.

To acquire a proper label, the sequences were presented to 15 human test persons. The sequences were presented as a video: after the play-back of a video the last image remained on the screen and the test person was asked to select a label. Thus, a fuzzy label for every sequence was created as the mean of the 15 different opinions. The resulting data collection showed to be highly imbalanced: the class “happiness” (105 samples) occurred four times more often than the class “fear” (25 samples) while “anger” (49 samples), “surprise” (91 samples), “disgust” and “sadness” (both 81 samples) are caught in between.

In our approach prominent facial regions such as the eyes, including the eye-brows, the mouth and for comparison the full facial region have been considered. For these four regions orientation histograms, principal components and optical flow features have been computed. Principal components are very well known in face recognition, and orientation histograms were successfully applied for the recognition of hand gestures [13]. To mitigate the impact of a subject’s individual facial form, the features based on orientation histograms and PCA, the individual face has to be eliminated, by subtracting a vector of the correct mean of the sequences locally. In order to extract the facial motion in these regions, optical flow features from pairs of consecutive images have been computed.

## 6 Experiments and results

In order to evaluate a team of individual classifiers, that has been trained based on unsupervised discovered categories, two competing approaches, using different category-systems has been implemented and evaluated. In general the implementations were designed in the following way: for both approaches, the individual classifiers were chosen to be HMM as they have proven to be eligible for the employed data in previous studies [14]. For all twelve available features-views on the data and for the considered classes (two classes in case of the fully supervised approach and number of clusters many in case of unsupervised processing step) a HMM was trained. Based on the probabilistic outputs of these models, a linear pseudo-inverse mapping to the true assigned label is computed. It is worth noting, that in case of the fully supervised approach this further mapping utilizes the very same labels as the latter one, whereas considering the approach using unsupervised learnt categories, these labels are new to the system.

Considering clustering, for every feature type the distances between all of the sequences were calculated using the HMM based distance measure as described in Section 3 and partitioned into ten clusters applying Ward’s method. Preliminary experiments have shown, that using ten partitions, the data can be divided in

well balanced clusters.

As in our experiments using fully supervised learning we considered solely two class discriminations, the output of individual classifier layer is a 24-dimensional vector. On the other hand, when using the unsupervised step, a probabilistic output for every of the ten categories is produced yielding a 120 dimensional probabilistic feature vector. In both cases, the target dimensionality of the fusion mapping is two as only pairings of two classes are considered.

<b>class. labels\ind. classifiers</b>	<b>unsupervised</b>	<b>supervised</b>
“hap.” vs. “anger”	<i>90.9</i>	84.4
“hap.” vs. “surprise”	<i>99.5</i>	96.9
“hap.” vs. “disgust”	<i>88.2</i>	81.2
“hap.” vs. “sadness”	<i>97.3</i>	90.3
“hap.” vs. “fear”	<i>74.6</i>	<i>80.8</i>
“anger” vs. “surprise”	<i>87.9</i>	85.7
“anger” vs. “disgust”	55.4	<i>64.6</i>
“anger” vs. “sadness”	70.0	<i>78.5</i>
“anger” vs. “fear”	<i>90.5</i>	64.9
“surprise” vs. “disgust”	<i>92.4</i>	82.0
“surprise” vs. “sadness”	<i>94.2</i>	87.2
“surprise” vs. “fear”	<i>66.4</i>	<i>76.7</i>
“disgust” vs. “sadness”	<i>75.9</i>	69.1
“disgust” vs. “fear”	<i>62.3</i>	63.2
“sadness” vs. “fear”	<i>87.7</i>	79.2

Table 1: Recognition rates of ten fold cross validation for the one-against-one classifiers in percent. The italic fonts indicate a line-wise maximum.

Table 1 shows the rate of correctly classified sequences for every combination of two classes and both proposed approaches. In ten of the 15 possible pairings of label systems, an improvement is observable, when the individual HMM is trained using the categories found by clustering. A closer examination of Table 1 shows, that considering classes “anger” and “fear” the partially unsupervised shows impaired classification rates.

Interpreting these findings, one could argue, that utilizing the data, which is considered as labelled” could have, in combination with the unsupervised preprocessing, beneficial effects for the following fusion and hence the over-all supervised classification in this approach. For the fully supervised classification, this data, which is not member of one of the selected classes, is obviously not available. On the other hand, positive effects are not observable in every case: especially when classes having fewer samples in the utilized data set are involved. This might be an indicator for in what ratio of labelled to unlabelled data brings benefit for a classifier and that this ration is not arbitrary.

## 7 Bottom line

In this study the usage of unlabelled data in a classification process is investigated. For this purpose, two classification approaches have been evaluated: one purely supervised and the other utilizing an unsupervised preprocessing step making use of all the available data. First results show, that doing so, improvements for the classifier system in terms of recognition rate can be produced. Future work will have to confirm these findings using different data sets.

## References

- [1] Lei Xu, Adam Krzyzak, and Ching Y. Suen. Methods of combining multiple classifiers and their applications in handwritten character recognition. *IEEE Transactions on Systems, Man and Cybernetics*, 22(3):418–435, 1992.
- [2] Ludmila I. Kuncheva. *Combining Pattern Classifiers: Methods and Algorithms*. Wiley-Interscience, 2004.
- [3] Lawrence R. Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. In *Proceedings of the IEEE*, pages 257–286, 1989.
- [4] Padhraic Smyth. Clustering sequences with hidden markov models. In *Advances in Neural Information Processing Systems*, volume 9, pages 648–654, 1997.
- [5] Manuele Bicego, Vittorio Murino, and Mario A.T. Figueiredo. Similarity-based clustering of sequences using hidden markov models. In *Machine Learning and Data Mining in Pattern Recognition, vol. LNAI 2734*, pages 86–95. Springer, 2003.
- [6] Andrew R. Webb. *Statistical Pattern Recognition*. John Wiley & Sons, Ltd, QinetiQ Ltd, Malvern, UK, 2 edition, 2003.
- [7] Gavin Brown and Ludmila I. Kuncheva. "good" and "bad" diversity in majority vote ensembles. In Neamat El Gayar, Josef Kittler, and Fabio Roli, editors, *MCS*, volume 5997 of *Lecture Notes in Computer Science*, pages 124–133. Springer, 2010.
- [8] David M. J. Tax, Martin van Breukelen, Robert P. W. Duin, and Josef Kittler. Combining multiple classifiers by averaging or by multiplying. *Pattern Recognition*, 33(9):1475 – 1485, 2000.
- [9] Ludmila Kuncheva, James C. Bezdek, and Robert P. W. Duin. Decision templates for multiple classifier fusion: an experimental comparison. *Pattern Recognition*, 34(2):299–314, 2001.
- [10] Friedhelm Schwenker, Christian Dietrich, Christian Thiel, and Günther Palm. Learning of decision fusion mappings for pattern recognition. *International Journal on Artificial Intelligence and Machine Learning (AIML)*, 6:17–21, 2006.
- [11] Roger Penrose. A generalized inverse for matrices. In *Proceedings of the Cambridge Philosophical Society*, volume 52, pages 406–413, 1955.
- [12] Takeo Kanade, Jeffrey Cohn, and Ying-Li Tian. Comprehensive database for facial expression analysis. In *Proceedings of the 4th IEEE International Conference on Automatic Face and Gesture Recognition (FG'00)*, pages 46–53, March 2000.
- [13] William T. Freeman and Michal Roth. Orientation histograms for hand gesture recognition. In *international workshop on automatic face and gesture recognition*, pages 296–301, 1994.
- [14] Miriam Schmidt, Martin Schels, and Friedhelm Schwenker. A hidden markov model based approach for facial expression recognition in image sequences. In Friedhelm Schwenker and Neamat El Gayar, editors, *4th IAPR TC3 Workshop on Artificial Neural Networks in Pattern Recognition (ANNPR 2010)*, LNAI 5998, pages 149–160. Springer, 2010.