

INCREMENTAL LEARNING IN HIERARCHICAL NEURAL NETWORKS FOR OBJECT RECOGNITION

Rebecca Fay, Friedhelm Schwenker and Günther Palm
University of Ulm
Department of Neural Information Processing
D-89069 Ulm, Germany

Keywords: Object recognition, hierarchical neural networks, incremental learning.

Abstract: Robots that perform non-trivial tasks in real-world environments are likely to encounter objects they have not seen before. Thus the ability to learn new objects is an essential skill for advanced mobile service robots. The model presented in this paper has the ability to learn new objects it is shown during run time. This improves the adaptability of the approach and thus enables the robot to adjust to new situations. The intention is to verify whether and how well hierarchical neural networks are suited for this purpose. The experiments conducted to answer this question showed that the proposed incremental learning approach is applicable for hierarchical neural networks and provides satisfactory classification results.

1 INTRODUCTION

When performing tasks in complex environments it is probable that a robot is confronted with objects it has not seen before and consequently cannot identify. This makes the ability to learn new objects during run time an important capability for advanced mobile service robots. The robot is then able to incrementally learn novel objects and thereby increase knowledge of its environment and adapt to new situations.

Without the ability to incrementally learn new objects the robot could only deal with predetermined objects learnt offline. All other objects could only be classified as unknown objects. In real world environments which are fairly complex and subject to numerous changes solely being capable of coping with previously learnt objects might not be sufficient.

The presented approach is an extension of our previous work (Knoblauch et al., 2004) which enables the model to deal with unknown objects by incrementally learning to categorise these objects.

2 HIERARCHICAL OBJECT RECOGNITION

The visual object recognition is implemented as a two stage process (Fay et al., 2004) where low resolu-

tion features such as colour and shape information are used to localise objects of potential interest and ensuing high resolution features such as edges, corners or T-junctions are extracted to analyse the part of the image containing the detected object in greater detail. These features are used in a trained neural network for object recognition. For the application at hand orientation histograms (Freeman and Roth, 1995) (Coppola et al., 1998) summing up all orientations (directions of edges represented by the gradient) within the region of interest are used as features for the object recognition.

2.1 Hierarchical Neural Networks

If neural networks are used for object recognition an object is represented by a number of features, which form a d dimensional feature vector x within the feature space $X \subseteq \mathbb{R}^d$. A classifier therefore realises a mapping from feature space X to a finite set of classes $C = \{1, 2, \dots, l\}$. A neural network is trained to perform a classification task using supervised learning algorithms. A set of training examples $S := \{(x^\mu, t^\mu), \mu = 1, \dots, M\}$ is presented to the network. The training set consists of M feature vectors $x^\mu \in \mathbb{R}^d$ each labelled with a class membership $t^\mu \in C$. During the training phase the network parameters are adapted to approximate this mapping as accurately as possible. In the classification phase un-

labelled data $x^\mu \in \mathbb{R}^d$ are presented to the trained network. The network output $K(x^\mu) = c \in C$ is interpreted as an estimation of the class corresponding to the input vector x .

The object classification is performed by a hierarchical neural network (Fay et al., 2004). It consists of several simple neural networks that are combined in a tree, i.e. the nodes within the hierarchy represent individual neural classifiers. For the application at hand RBF networks were used as neural classifiers. The basic idea of hierarchical neural networks is the hierarchical decomposition of a complex classification problem into several less complex ones. This yields hierarchical class grouping whereby the decision process is split into multiple steps exploiting rough to detailed classification. The hierarchy emerges from recursive partitioning of the original set of classes C into several disjoint subsets C_i until subsets consisting of single classes result. C_i is the subset of classes to be classified by node i , where i is a recursively composed index reflecting the path from the root node to node i . The subset C_i of node i is decomposed into s_i disjoint subsets $C_{i,j}$, where $C_{i,j} \subset C_i$, $C_i = \cup_{j=0}^{s_i-1} C_{i,j}$ and $C_{i,j} \cap C_{i,k} = \emptyset$. The total set of classes $C_0 = C$ is assigned to the root node. Consequently nodes at higher levels of the hierarchy typically classify between larger subsets of classes whereas nodes at the lowest level discriminate between single classes. This divide-and-conquer strategy yields several simple classifiers, that are more easily manageable than one extensive classifier. These simple classifiers can be amended much more easily to the decomposed simple classification tasks than one classifier could be adapted to the original complex classification task. Furthermore different feature types X_i are used within the hierarchy. For each classification task the feature type that allows for the best discrimination is chosen. An example of such a hierarchy is shown in figure 1. The feature types X_i and X_j used for node i and j can be identical, i.e. both can use the same feature type.

2.2 Hierarchy Generation

The hierarchy is generated by unsupervised k-means clustering. In order to decompose the set of classes C_i assigned to one node i into s_i disjoint subsets a k-means clustering is performed with all data points $\{x^\mu \in X_i | t^\mu \in C_i\}$ belonging to these classes. Depending on the distribution of the classes across the k-means clusters s_i disjoint subsets $C_{i,j}$ are formed. One successor node j corresponds to each subset. For each successor node j again a k-means clustering is performed to further decompose the corresponding subset $C_{i,j}$. The k-means clustering is performed for each feature type. The different clusterings are eval-

uated and the clusterings which group data according to their class labels are preferred. Since the k-means algorithm depends on the initialisation of the clusters, k-means clustering is performed several times per feature type. In this study the number of k-means clustering per feature type was 10.

The number of clusters k must be at least the number of successor nodes or the number of subsets s respectively but can also exceed this number. If the number of clusters is higher than the number of successor nodes, several clusters are grouped together so that the number of groups equals the number of successor nodes. For reasons of simplicity in the following only the case where the number of clusters k equals the number of successor nodes s is considered. The valuation function prefers clusterings that group data according to their class labels. Clusterings where data are uniformly distributed across clusters notwithstanding their class labels receive low ratings. Furthermore clusterings are preferred which evenly divide the classes. Thus the valuation function rewards unambiguity regarding the class affiliation of the data assigned to a prototype as well as uniform distribution regarding the number of data points assigned to each prototype.

To generate the hierarchy at first the set of all classes is assigned to the root node. Starting with a clustering on the complete data set the set of classes is divided into subsets. Each subset is assigned to a successor node of the root node. Now the decomposition of the subsets is continued until no further decomposition is possible or until the decomposition does not lead to a new division. An example of a classification hierarchy is shown in figure 1.

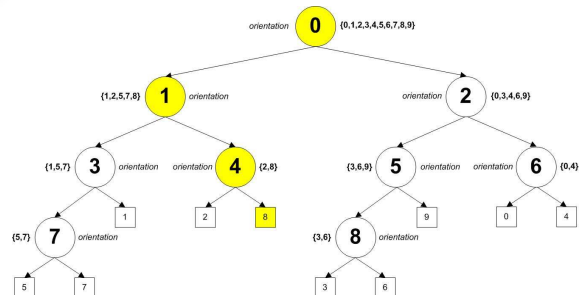


Figure 1: Classifier hierarchy generated for the classification of 10 classes of the COIL20 data set using orientation histograms as feature type. Each node within the hierarchy represents a neural network which is used as a classifier. The end nodes represent classes. To each node a feature type and a set of classes is assigned. The corresponding neural network uses the assigned feature type to discriminate between the assigned classes. The highlighted path shows the nodes activated during the classification of a sample that is classified as class 8.

2.3 Training and Classification

The hierarchy is trained by separately training the individual classifiers with the data $\{x^\mu \in X_i | t^\mu \in C_i\}$ that belong to the subsets of classes assigned to each classifier. For the training the respective feature type X_i identified during the hierarchy generation phase is used. The data will be relabelled so that all data points of the classes belonging to one subset $C_{i,j}$ have the same label j , i.e. $\tilde{t}^\mu = j, x^\mu \in X_i, t^\mu \in C_{i,j}$. The number of input nodes of the single classifiers is defined by the dimension d_i of the respective feature type X_i assigned to the corresponding node i . The number of output nodes equals the number of successor nodes s_i . The classifiers are trained using supervised learning algorithms. The classifiers within the hierarchy can be trained independently, i.e. all classifiers can be trained in parallel.

Within the hierarchy different types of classifiers can be used. Examples of classifiers would be radial basis function (RBF) networks, linear vector quantisation classifiers (Simon et al., 2002) or support vector machines (Schwenker, 2001). We chose RBF networks as classifiers. They were trained with a three phase learning algorithm (Schwenker et al., 2001).

The classification result is obtained similar to the retrieval process in a decision tree (Duda et al., 2001). Starting with the root node the respective feature vector of the object to be classified is presented to the trained classifier. By means of the classification output the next classifier to categorise the data point is determined, i.e. the classifier j^* corresponding to the highest output value $o(j^*)$ is chosen such that $j^* = \operatorname{argmax}_{j=1..s_i}(o(j))$. Thus a path through the hierarchy from the root node to an end node is obtained which not only represents the class of the object but also the subsets of classes to which the object belongs. Hence the data point is not presented to all classifiers within the hierarchy and the hierarchical decomposition of the classification problem yields additional intermediate information.

If only intermediate results are of interest it might not be necessary to evaluate the complete path. In order to solve a task it might be sufficient to know whether the object to be recognised belongs to a set of classes and the knowledge of the specific category of the object might not add any value. If the task for example is to grasp a cup, it is not necessary to distinguish between red and green cups. Moreover, when looking for a specific object it might in some cases not be necessary to retrieve the final classification result if a decision at a higher level of the hierarchy already excludes this object.

2.4 Incremental Learning of New Classes

While performing tasks a robot is likely to encounter unfamiliar objects. Hence the ability to learn new objects plays an important role. Given the situation that an object is in the robot's visual field and it is told the name of the new object by an instructor, the robot might look at the object from different points of view and thereby generate a few data samples $\tilde{S} := \{(x^\nu, \tilde{c}), \nu = 1, \dots, N\}$ of the new class $\tilde{c} \notin C$. So compared to the already known classes C the number of samples N of the new class \tilde{c} is appreciably lower. This scenario forms the basis for the incremental learning approach developed here. In the overall system the online learning is triggered by instruction sentences starting with "this is" followed by an object name.

In order to quickly achieve results the learning is performed in two stages. In the first stage fast but less sophisticated methods are used to obtain initial results, i.e. the novel objects are learnt but the recognition rate might be weak. In this first stage the recognition rate can be improved by using a similar method to retrain the new object with additional data gained by taking different views of the object. In a second stage more complex algorithms are used to adapt the system and to further improve the classification results. This retraining is more advanced but also more time consuming.

At first it is necessary to identify whether the presented object is already known or not. This is accomplished by presenting the new data to the trained classifier and taking the strength of the classifier response into account. Thereby a strong response is considered as an unambiguous decision and weak responses indicate a dubious decision which could be evoked by unknown classes or if the object to classify bears resemblance to more than one class. The thresholds for this are derived from the classifier responses when testing the known data. If the new data is unambiguously classified as one class c without exception it is assumed that the object is already known and belongs to class c . Otherwise the object is regarded as a hitherto unidentified object.

If an object is identified as unfamiliar it is learnt by fitting it into the hierarchy and if necessary retraining the affected nodes. The new class $\tilde{c} \notin C$ associated with the unknown object is inserted as a new leaf. The position of the new leaf is determined by classifying all new data and evaluating the corresponding paths through the hierarchy. The leaf will be inserted where the paths start to diverge. As complete identicalness for all data cannot be presumed even at the root node since the network has not been trained with this data a certain variance needs to be considered. Otherwise

the new leaf would in most instances be added at the root node. Therefore the common path is successively determined similar to retrieving the classification result. Starting with the root node all new samples are presented to the corresponding classifier. Depending on the classification results either the next node in the path is determined or the search is stopped and the new class is added as a new leaf at this node. If all new samples x^v are assigned to one successor node j by classifier i node i is added to the common path without retraining and classifier j is the next classifier. If not all but a significant majority of the data points is assigned to the same successor node j this classifier is retrained using the samples of the known objects as well as the new samples and node i is then added to the common path. If there is no clear decision or if all successor nodes of node i are end nodes the new class is added as an additional end node to this node and the node is afterwards retrained. Because of the way the hierarchy is built only part of the hierarchical classifier needs to be amended. The rest of the hierarchy remains unchanged. Figure 2 depicts how a new class is inserted into the hierarchy by means of the proposed incremental learning approach.

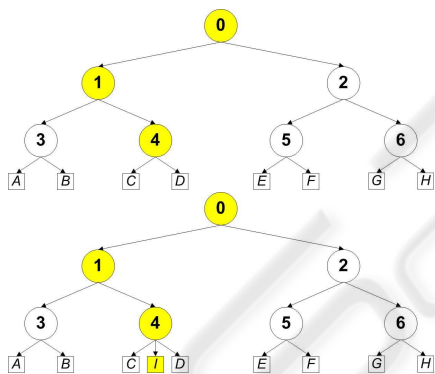


Figure 2: Example of the incremental learning of a new class I. At first the position where to insert the new class is determined. Then the new class is added as a new leaf and the affected nodes are retrained. Here the new class is added to node 4. The classifiers 0 and 1 are retrained with the samples of the additional class I. One additional node is added to the output layer of classifier 4 and the classifier is then retrained.

The retraining or incremental training of the classifiers is conducted by adding a new neuron to the hidden layer and then retraining the output weights with the joint sample set of old and new samples. The centre of the new prototype is determined by the mean of all new samples. The width of the corresponding gaussian function is set to the mean distance of the new samples to the centre and the new output weights are learnt by calculating the pseudo-inverse (Penrose, 1955).

A similar mechanism is applied when retraining already learnt classes. The only differences are that no additional leaf is added and that the path through the hierarchy is already known. The single classifiers on this path are retrained if there are any incorrect or ambiguous decisions.

Scenarios for retraining an already learnt class could be when a class is only represented by few samples or the classification performance for this class is not satisfactory. The classification rate of novel classes is likely to be lower than the classification rate for the previously learnt classes. Thus a retraining of all affected nodes can yield improved classification results once sufficient additional samples for the novel class are available. Another reason could be that a new instance of a class with noticeably different characteristics has to be learnt.

The online learning phase is followed by an offline learning phase where more sophisticated learning algorithms such as three phase learning are used which will further improve the classification performance. All nodes to which the novel classes have been assigned are newly trained. These algorithms would be too time consuming for usage during run time.

This incremental learning approach can also be used for non-hierarchical neural networks. However, here it is not necessary to determine the position where to insert the new leaf as non-hierarchical networks can be regarded as hierarchical networks consisting only of one node, but a retraining of the complete network takes place whereas only parts of the hierarchical network is amended.

3 RESULTS

By means of classification experiments the suitability of hierarchical neural networks for extension was examined. It should be verified whether new classes that are only represented by a few samples can be learnt sufficiently well in moderate time and whether it is possible to learn new classes without negatively affecting the classification performance of already learnt classes.

The incremental learning approach has been evaluated using the Columbia Object Image Library (COIL20) (Nene et al., 1996) consisting of 20 objects. Therefore 10-times 10-fold cross-validation experiments have been conducted on the camera images of the COIL20 image library.

Another approach for learning new objects would be a complete redesign of the classifier. This means that everything learnt so far is dropped and the classifier is rebuilt and trained with the joined samples of the known and unknown objects. For the classifier training only simple training algorithms such as two-

phase learning (Schwenker et al., 2001) can be used as sophisticated learning algorithms would be too time-consuming and hence not applicable for online learning. Compared to the incremental learning approach this method features longer training time as a complete rebuild of the hierarchy and training of all nodes within the hierarchy are required. The approach proposed has been compared to this method.

In order to facilitate comparability of the results of the different approaches for all experiments the same division of the data for the cross-validation has been used.

For the experiments 10 classes of the 20 classes of the COIL20 data set were chosen to represent the familiar objects. The remaining 10 classes formed a pool of potentially new objects. To test the incremental learning approach proposed the first 10 classes were used to generate and train a hierarchy. To train the hierarchy sophisticated learning algorithms were used. Here the three-phase learning for RBF networks was used. The hierarchy structure was held constant for reasons of comparability while the hierarchy training was subject to cross-validation. This means the same hierarchy was trained 10×10 times resulting in 100 trained hierarchies. These trained hierarchies for the classification of the first 10 classes formed the basis for the incremental learning experiments. Each of the 10 new classes was then learnt incremental so that $10 \times 10 \times 10$ hierarchies for the classification of 11 classes emerged. For the incremental learning 10 samples of the unknown class were used compared to 64 to 65 samples (depending on the specific cross-validation run) of the other classes. In the test data set all classes are represented with the same number of samples.

With the same partition of the samples the alternative method which completely rebuilds the hierarchy was evaluated. Therefore 10 hierarchies for the classification of 11 classes each were generated. For each of the 10 hierarchies 10×10 cross validation experiments were conducted.

Both approaches showed essentially the same classification quality. The incremental learning approach had an average classification rate of $95.94 \pm 2.96\%$ and the average classification rate of the alternative method was $95.54 \pm 2.56\%$. Despite being less extensive the incremental learning approach achieved the same results. Figure 3 visualises these results as box plots.

The confusion matrix for the incremental learning experiments displayed in figure 4 shows that although being represented by a significantly lower number of samples the classification rates of the new classes is equal to the classification rates of the primarily learnt classes.

Having a look at the positions within the hierarchy were the new classes were added, it could be found

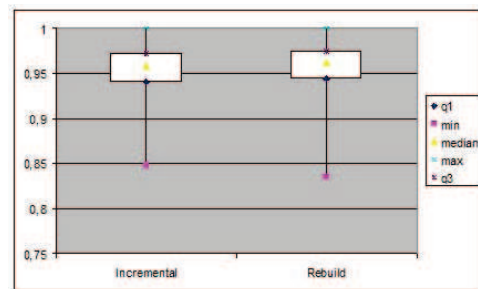


Figure 3: Box plot of the classification results of the two alternative approaches for learning new classes. The classification results do not differ significantly.

	0	1	2	3	4	5	6	7	8	9	new
0	0.8742	0	0	0.278	0	0	0	0	0	0.014	12.29
1	0	0.9456	0	0.111	0	0	0	0	0	0	5.333
2	0	0	0.9683	0.083	0	0	0	0	0	0	3.083
3	0	0.014	0	0.9706	0	0	0	0	0	0	2.931
4	0	0	0	0	1.00	0	0	0	0	0	0
5	0	0.306	0	0.097	0	0.9749	0	0	0	0	2.403
6	0	0	0	0.347	0	0	0.9899	0	0	0	0.667
7	0	0	0	0	0	0	0	0.9986	0	0	0.139
8	0	0	0	0	0	0	0	0	0.9517	0	4.833
9	0	0	0	0.083	0	0	0	0	0	0.9193	7.986
new	0	0.917	0	0.667	0	0	0	0	0.667	1.375	96.38

Figure 4: Confusion matrix for the experiments utilising incremental learning.

that in the most instances a class was mainly assigned to one node. If this is not the case then at least the class was added to nodes lying on the same path of the hierarchy as can be observed for class 13. The only exception to this is class 11, which is distributed over different pathes of the hierarchy. It can also be observed that in the majority of cases the classes were added to nodes in deeper levels of the hierarchy. No classes were added to the root node. Figure 5 illustrates for each new class how often it was added to which node of the hierarchy shown in figure 1.

These results show that the new classes are not added arbitrarily to the hierarchy but for each class preferred positions emerge. The noticeable frequency of node 8 is likely to be a characteristic of the COIL20 data set.

4 RELATED WORK

An example for other incremental learning approaches are ART networks. ART networks (Grossberg, 2000) (Carpenter and Grossberg, 2002) allow for online learning of evolving data sets. If a presented sample is similar enough to an learnt prototype this prototype is adjusted to the sample, otherwise a new prototype is defined by the sample. In the ARAM model (Tan, 1995) new classes can be learnt

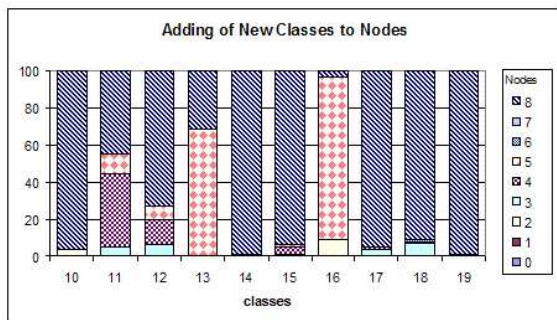


Figure 5: Adding of new classes to nodes. For each of the 10×10 cross-validation experiments conducted for each of the new classes 10 to 19 a leaf was added to different nodes of the classification hierarchy. For each new class is shown in percent to which node the corresponding leaf was added.

while preserving previously learnt classes and so the stability-plasticity dilemma is regarded. However, ART networks are non-hierarchical networks and they consider new samples one at a time.

5 CONCLUSIONS

The proposed approach has proven functional. Although the networks were trained with only a few samples of the new classes they were able to classify the new class and no considerable deterioration of the classification results of the former classes could be observed.

The experiments conducted showed encouraging results. It could be proved that hierarchical neural networks are suitable for incremental adaption. New objects can be learnt with sufficient classification rates and in adequate time.

The proposed approach enables the robot to deal with varying object categories in addition to the pre-defined categories. New objects can be learnt rather quickly with satisfactory quality. The quality can even be increased by retraining with additional samples and in an offline phase sophisticated learning algorithms are used to further improve the classification quality.

ACKNOWLEDGEMENT

This research was supported in part by the European Union award #IST-2001-35282 of the MirrorBot project.

REFERENCES

- Carpenter, G. A. and Grossberg, S. (2002). Adaptive Resonance Theory. In Arbib, M., editor, *The Handbook of Brain Theory and Neural Networks*, pages 87–90. MIT Press, Cambridge, 2nd edition.
- Coppola, D. M., Purves, H. R., McCoy, A. N., and Purves, D. (1998). The distribution of oriented contours in the real world. *Proceedings of the National Academy of Sciences USA*, 95(7):4002–4006.
- Duda, R. O., Hart, P. E., and Stork, D. G. (2001). *Pattern Classification*. John Wiley & Sons, New York, 2nd edition.
- Fay, R., Kaufmann, U., Schwenker, F., and Palm, G. (2004). Learning object recognition in a neurobotic system. *3rd IWK Workshop SOAVE2004 - SelfOrganization of Adaptive behavior, Illmenau, Germany*, pages 198–209.
- Freeman, W. T. and Roth, M. (1995). Orientation histograms for hand gesture recognition. In *IEEE International Workshop on Automatic Face- and Gesture-Recognition*, pages 296–301, Zürich, Switzerland.
- Grossberg, S. (2000). Adaptive resonance theory. Technical Report TR-2000-024, Center for Adaptive Systems and Department of Cognitive and Neural Science, Boston University.
- Knoblauch, A., Fay, R., Kaufmann, U., Markert, H., and Palm, G. (2004). Associating Words to Visually Recognized Objects. In Coradeschi, S. and Saffiotti, A., editors, *Anchoring symbols to sensor data. Papers from the AAAI Workshop. Technical Report WS-04-03*, pages 10–16. AAAI Press, Menlo Park, California.
- Nene, S. A., Nayar, S. K., and Murase, H. (1996). Columbia object image library (coil-20). Technical Report Technical Report CUCS-005-96, Columbia University.
- Penrose, R. (1955). A generalized inverse for matrices. *Mathematical Proceedings of the Cambridge Philosophical Society*, 51:406–413.
- Schwenker, F. (2001). Solving multi-class pattern recognition problems with tree structured support vector machines. In Radig, B. and Florczyk, S., editors, *Mustererkennung 2001*, pages 283–290. Springer.
- Schwenker, F., Kestler, H. A., and Palm, G. (2001). Three learning phases for radial-basis-function networks. *Neural Networks*, 14:439–458.
- Simon, S., Schwenker, F., Kestler, H. A., Kraetzschmar, G. K., and Palm, G. (2002). Hierarchical object classification for autonomous mobile robots. In *International Conference on Artificial Neural Networks (ICANN)*, pages 831–836.
- Tan, A.-H. (1995). Adaptive resonance associative map. *Neural Networks*, 8(3):437–446.