

A QUANTITATIVE STUDY OF EXPERIMENTAL NEURAL NETWORK LEARNING ALGORITHM EVALUATION PRACTICES

Lutz Prechelt (prechelt@ira.uka.de)

Universität Karlsruhe, Germany

Abstract. *113 articles about neural network learning algorithms published in 1993 and 1994 are examined for the amount of experimental evaluation they contain. Every third of them does employ not even a single realistic or real learning problem. Only 6% of all articles present results for more than one problem using real world data. Furthermore, one third of all articles does not present any quantitative comparison with a previously known algorithm. These results indicate that the quality of research in the area of neural network learning algorithms needs improvement. The publication standards should be raised and easily accessible collections of example problems be built.*

INTRODUCTION

A large body of research in artificial neural networks is concerned with finding good learning algorithms to solve practical application problems. Such work tries to improve for instance the quality of found solutions (generalization), the probability of convergence, the ease of use, the learning speed, or some combination thereof. Currently, there exists no theory that quantitatively predicts the behavior of a new algorithm compared to other algorithms for any of these criteria. Consequently, experimental evaluation¹ is needed to validate any claims of improvement made for a new algorithm or to characterize under which circumstances improvements can be expected.

It seems that such evaluation is often not performed thoroughly enough, even in articles published by leading journals. Motivated by this impression, I decided to investigate this hypothesis by studying the current research practice empirically. In a recent study by Tichy et al. (1) about experimental evaluation in computer science publications, the journal *Neural Computation* produced quite good results, far above average. However, the only measure used in that work was the fraction of article space devoted to the evaluation and the articles considered were not only those about learning algorithms. The approach taken in the present study is more concrete at assessing the quality of an evaluation. I review the set of all articles presenting learning algorithms

for practical problems that appeared in two renowned neural network journals in 1993 and the first half of 1994. In each article, the number of problems used in the algorithm evaluation and the number of other algorithms used for comparison were counted. While high numbers resulting from such counting cannot prove that the evaluation has high quality, low numbers prove that the quality is low.

The articles under consideration are from the two oldest journals dedicated to neural network research, namely *Neural Networks* (NN), the official journal of the International Neural Network Society, published by Elsevier, and *Neural Computation* (NC), published by MIT Press. From *Neural Networks*, all articles of volume 6 (1993) and all articles from numbers 1 to 5 of volume 7 (1994) were used. From *Neural Computation*, all articles of volume 5 (1993) and all articles from numbers 1 to 4 of volume 6 (1994) were used.

The subsequent sections present the methodology and limitations of the study, the obtained results, and the conclusions drawn.

METHODOLOGY

Approach

The objective of the present study is to determine the quality of current algorithm evaluations. As a measure of quality we use the number of problems and compared algorithms used in an evaluation. The exact criteria are described in the next section. We consider the quality of the evaluation to be low when these numbers are low. If the numbers are high, no statement of quality can be made with this method. The rationale of this approach is to make the results as objective and reproducible as possible.

Method

1. Each article was classified into one of the following categories.

Theory. Articles belong to the “Theory” category if and only if the major contributions made by the paper are formally proven propositions.

¹In this report, I will use the term *evaluation* to mean *experimental evaluation*.

Modeling. Articles predominantly concerned with the formal modeling of some aspects of natural neural networks, or with discussing the properties of such models, or with other aspects of biological plausibility belong to the “Modeling” category.

Algorithm. Articles whose main contribution is the design of a new learning algorithm to be applied to practical problems form the “Algorithm” category². Empirical studies comparing several known algorithms and application papers presenting architectures for applying known algorithms to a particular problem field are also included here, since they are quite rare (only 5% of the category).

Other. All articles that do not fit into any of the above categories are put into the “Other” category. This includes surveys and papers on electronic neural network hardware.

“If in doubt, leave it out.” In all borderline cases, papers were *not* classified as Algorithm in order to avoid a negative bias in the data due to papers that were not meant to make an algorithm contribution and, thus, lack proper evaluation. In particular, any paper appearing in Neural Computation that was marked to be a “Note” and that would have been an Algorithm paper by its topic was classified as Other in order to avoid a negative bias in the data due to papers that were simply too short to contain proper evaluation.

2. After the category of each article was determined, only the articles from the Algorithm category were used in the study. Each Algorithm article was reviewed to determine the two key metrics used in the study, namely

- the number of different learning problems (data sets) used in the evaluation and
- the number of known algorithms a proposed algorithm is compared to.

For a more meaningful discussion, each learning problem is classified to be either an artificial, a realistic, or a real problem.

Artificial problems are those whose data is generated synthetically based on some simple logic or arithmetic formula.

Realistic problems also consist of synthetic data, but are generated by a model with properties similar to what can be found in real problems. Only the following three types of data generation procedures yield what is considered realistic problems: firstly, data generation using a complex and realistic mathematical model of a physical system such as a cart/pole system or robot kinematics; secondly, data

generation by chaotic mathematical processes, such as the Mackey-Glass equation; and thirdly, data generation by stochastic processes, such as mixtures of Gaussian random variables.

Realistic problems are useful to assess the behavior of an algorithm on problems with known properties; they provide the best way to *characterize* the kinds of problems for which an algorithm will yield good results.

Real problems consist of data that represents actual observations of phenomena in the physical world. Such data tends to contain some amount of errors and noise. Most importantly and in contrast to realistic artificial data, real data usually has characteristics that are not completely known (surprising features). We want learning algorithms to cope well with problems whose characteristics are partially unknown; how well they do can best be tested with real data.

Synthetic variations of the same problem count as a separate problem only if it is plausible to expect that two algorithms may compare very different on the variation than on the original problem. In many cases, two variations of a problem were found: one with and one without noise in the data. A very different problem representation is another kind of problem variation that counts as a separate problem. What exactly “very different” means cannot be quantified, but I did my best to apply constant criteria throughout the study.

To *use* a problem in an evaluation means to report any kind of quantitative data about the behavior of the proposed algorithm on this problem, for instance learning speed, convergence probability, training set error, or test set error.

The algorithms used for comparison were originally discriminated to be either neural network algorithm or other algorithms. Since that discrimination is fuzzy, however, the separation is dropped in the discussion of the results. The count includes all algorithms not introduced in the article in question; algorithms that are newly proposed in an article are not counted. Articles presenting comparative empirical studies of known algorithms had all algorithms counted. When an article introduces several new algorithms at once, all algorithms used for a comparison to *any* of the new ones are counted, i.e., an algorithm used for comparison is counted even if it is not compared to all of the new algorithms.

Limitations

The method described above does not allow for a quantitative judgement of the overall quality of an evaluation. Even if many problems and compared algorithms are used, the relevance of the results may

²The word Algorithm, with capital A, will be used throughout this report to refer to the category.

still be low due to irrelevant performance measures, irrelevant or biased problems, improper description of the setup, or other methodological errors. The assumption used in the approach is *not* that a large number of problems and compared algorithms in an article implies high evaluation quality, but instead that a small number implies low evaluation quality. The counting criteria themselves are biased towards finding large numbers.

An absolute quality measure is not required, since all this study is meant to do is investigating the hypothesis that algorithm evaluations are often of low quality. No attempt will be made to quantify what low quality means, because any such quantification would necessarily be arbitrary. Instead, we will reject the hypothesis unless we find subjectively overwhelming evidence for it. Hence, the approach of the study is quite conservative.

Nevertheless, a few remarks must be made on possible objections against the approach.

1. *An algorithm proposed for a narrow application domain does not allow for a wide variety of test problems.* This is true, but is not the issue debated here. Even for a very specialized algorithm, a number of different incarnations of problems from its domain can be found and should be investigated. For instance, variations of a problem obtained by significantly changing a major parameter such as the resolution of the data would be counted as separate problems. Only the number of problems is judged, not their variety.
2. *Often no algorithms can be found to be compared to an algorithm proposed for a narrow application domain.* Maybe no other specialized algorithms can be found. But it is nevertheless interesting to see how much improvement the new algorithm represents compared to known general purpose algorithms. Thus, such algorithms should be used for comparison.
3. *Algorithms solving a problem for which no solution was previously known cannot be compared to others.* This is true, but it hardly ever applies; I did not observe any instance of such an algorithm in the whole sample investigated in this study, although arguably there are a few borderline cases.
4. *Totally new approaches to a problem do not allow for comparison.* Why not? If the approach was made for its assumed utility, a comparison is the best means to assess it. Otherwise the article should not claim utility and would then be classified as Modeling in this study.
5. *Often a thorough evaluation is simply too much work.* The result of scientific work should be knowledge. An algorithm about whose behavior too little knowledge is available is no proper

scientific contribution. Experimental evaluation may be a lot of work, but it needs to be done.

6. *I believe that your data contains many errors.* Probably there is a considerable number of errors in my data. See (1) for a discussion and estimation of the precision to be expected from a study like the present one. However, as we will see below, the conclusions from this study do not change even if a large margin of error is assumed.

RESULTS AND DISCUSSION

The raw data obtained during the study is presented in Prechelt (2). In this section, I will present only the most prominent findings. Since the differences between Neural Networks and Neural Computation are quite small in most respects covered here, I will discuss the set of all Algorithm articles studied as a whole.

Let us first have a look at the total number of problems used in the evaluation. This is depicted in figure 1.

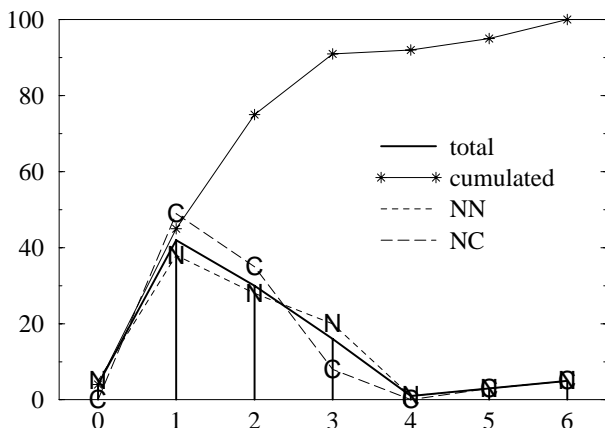


Figure 1: Percentage y of Algorithm articles that use a total of x different problems for the evaluation.

The figure is to be read as follows. On the abscissa (x -axis), we find the article classes from “0 problems used” up to “5 problems used”. The last point, $x = 6$, stands for “6 or more problems used”. The ordinate value (y -value) indicates the percentage of articles belonging to the class. The curve drawn as a thick line indicates the value for the total of all Algorithm articles found, while the dashed lines show the corresponding data for Neural Networks (NN) and Neural Computation (NC), respectively, alone. The starred line is the accumulation of the values on the thick line from left to right; it can be used to read quantiles. All other figures have the same structure.

As we see, 4% of all articles do not have any experimental evaluation and only 25% use more than

two problems for the evaluation. While it is surprising enough that any Algorithm article without experimental evaluation can be published in a renown journal, it is even more staggering how few articles use a broad set of problems. Only 9% of all articles use more than three problems.

Now let us differentiate this data by problems being either artificial, realistic, or real as defined above. Figure 2 shows the number of artificial problems used. No special remark is to be made here, since

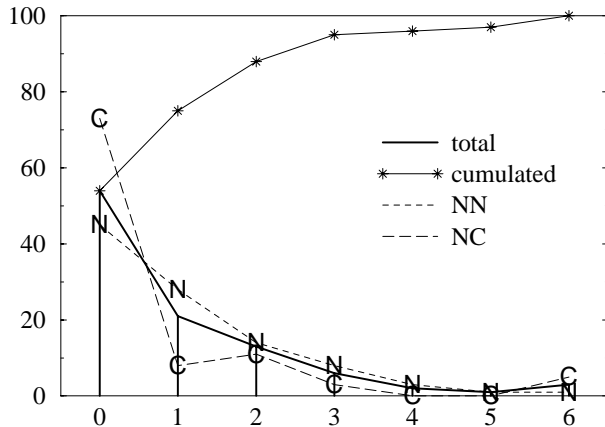


Figure 2: Percentage y of Algorithm articles that use x different artificial problems for the evaluation.

artificial problems should only serve for the illustration (as opposed to the evaluation) of an algorithm; a large number of artificial problems in an article is neither good nor bad. 20 articles (18%) employed the “grandfather” of all neural network problems, the XOR or n -bit parity.

Figure 3 shows the number of realistic problems used per article. As mentioned before, such problems are

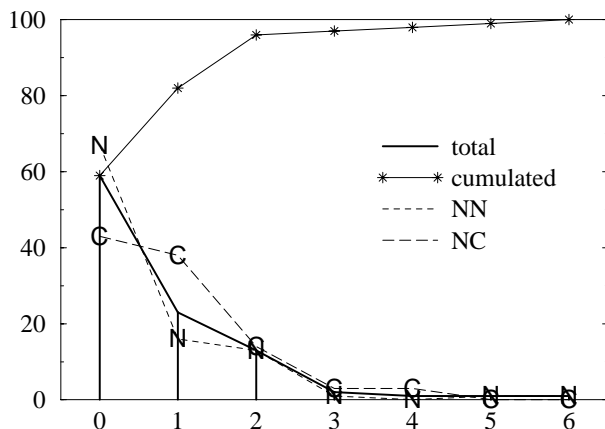


Figure 3: Percentage y of Algorithm articles that use x different realistic problems for the evaluation.

useful to explore an algorithm on data whose properties are realistic, yet exactly known. Despite that usefulness, 59% of all articles do not use any realistic

problem, only 4% use more than two, and 3% more than three. As we see, an experimental exploration of the question “For which kinds of problems is this algorithm best suited?” is hardly ever done.

Figure 4 shows the number of real problems used per article. Of course, nobody can say how re-

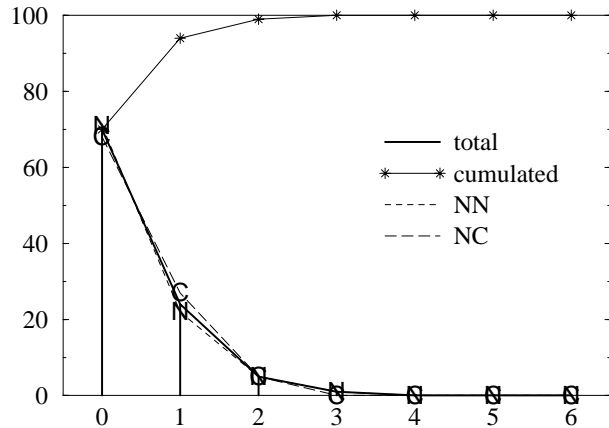


Figure 4: Percentage y of Algorithm articles that use x different real problems for the evaluation.

sults on one real problem (or, for that matter, 15 real problems) generalize to other problems, but it is also impossible to say exactly how the performance on realistic problems will generalize to real problems. Thus, it should at least be verified that an algorithm performs well for *some* real problems, as real problems are the only tests of a learning algorithm that are *guaranteed* to have at least some practical relevance (namely for the exact problem tested). Another reason is that real data tends to have some totally unexpected features that artificially generated data, even if otherwise realistic, lacks. However, the use of real problems in the articles of the study is deplorably rare. 70% of all articles do not use any real problem, only 1% use more than two, and not a single one was found using more than three.

Even when summing the number of realistic and real problems used in each article, as depicted in figure 5, a huge fraction of all articles is devoid of a reasonable number of test problems. 34% of all articles use zero realistic *and* zero real problems, 6% use more than two and a mere 3% use more than three.

The situation does not look much better when one considers the number of other algorithms used for comparison, as shown in figure 6. As much as 34% of all articles feature no comparison with other algorithms at all; only 19% compare to more than two known algorithms. This would not be a problem if everybody used standardized problems in standardized setups, but for the realistic and real problems this is not the case — it is very rare today that two different articles publish directly comparable results for the same problem. Without such comparability,

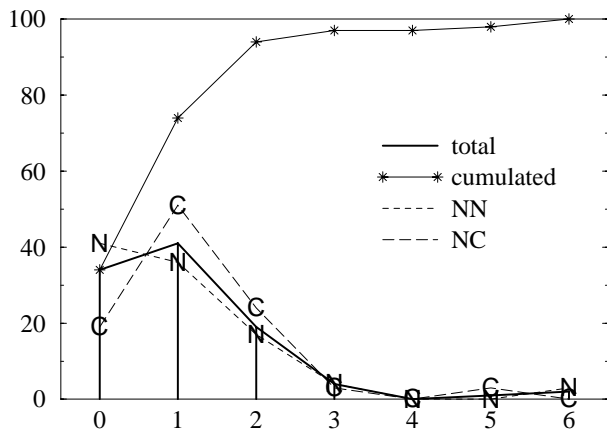


Figure 5: Percentage y of Algorithm articles that use x different realistic or real problems for the evaluation.

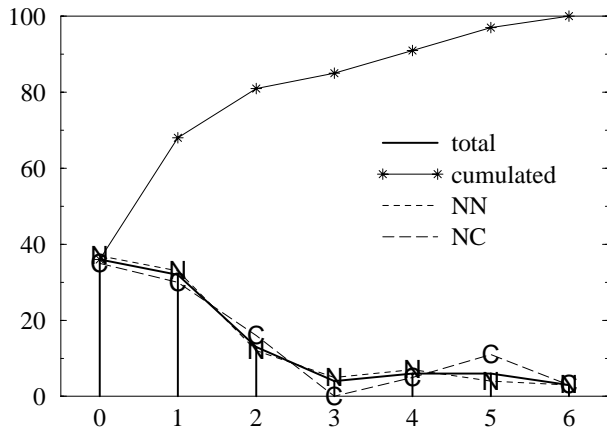


Figure 6: Percentage y of Algorithm articles that use x different known algorithms for comparison.

however, the above number means that for one out of every three articles the evaluation performed would better be called a naval inspection.

CONCLUSION

Assume that we set the following very modest standard. *An algorithm evaluation is called acceptable if it uses a minimum of two real or realistic problems and compares the results to those of at least one alternative algorithm.*

Then as much as 85% of Algorithm articles published in NN and NC do *not* meet this standard!

This result indicates that today new neural network learning algorithms are often published in a form that does not represent useful and validated knowledge. These articles present an idea of the kind “This is a way to tackle certain learning problems.”, but they do not tell us what we have to expect if we really try that idea. Instead, each article presenting a new algorithm should give at least a preliminary answer to the questions “For what kinds of problems

does the new algorithm work well or not well?” and “Under what conditions should we prefer the new algorithm over previously known ones?”. This information is essential if the publication of the algorithm is meant to be a scientific progress.

I believe the following steps should be taken to improve on the current situation.

1. Editors and reviewers set significantly higher standards for the experimental evaluation of a new learning algorithm. Articles that do not meet these standards are usually rejected.
2. Researchers reserve sufficient resources for a thorough experimental evaluation of their algorithms.
3. The research community prepares and uses public collections of example problems from all relevant fields in order to simplify algorithm evaluations. Re-use of example problems is also a prerequisite for broad comparisons of algorithms. Only a few fields such as speech recognition and optical character recognition do already have such collections.
4. Standard experimental setups and result presentation formats are developed to improve comparability and reproducibility of evaluation results.

Without these improvements, progress in the learning algorithm field will be significantly slower than it could be.

1. Tichy W.F., Lukowicz P., Prechelt L. and Heinz E.A., 1995, “Experimental evaluation in computer science: A quantitative study”, *Journal of Systems and Software*, 9–18. Also as Technical Report 17/94, Fakultät für Informatik, Universität Karlsruhe, Germany, August 1994, anonymous ftp: /pub/papers/techreports/1994/1994-17.ps.Z on ftp.ira.uka.de.

2. Prechelt L., 1994, “A study of experimental evaluations of neural network learning algorithms: Current research practice”, Technical Report 19/94, Fakultät für Informatik, Universität Karlsruhe, Germany, August 1994, anonymous ftp: /pub/papers/techreports/1994/1994-19.ps.Z on ftp.ira.uka.de.