

Neural learning of approximate simple regular languages

Mikel L. Forcada and Antonio M. Corbí-Bellot,
*Departament de Llenguatges i Sistemes Informàtics,
Universitat d'Alacant, E-03071 Alacant, Spain.*

and

Marco Gori and Marco Maggini,
*Dipartimento di Ingegneria dell'Informazione,
Università degli Studi di Siena,
via Roma, 56, I-53100 Siena, Italy.*

Abstract. Discrete-time recurrent neural networks (DTRNN) have been used to infer DFA from sets of examples and counterexamples; however, discrete algorithmic methods are much better at this task and clearly outperform DTRNN in space and time complexity. We show, however, how DTRNN may be used to learn not the *exact* language that explains the whole learning set but an *approximate* and much simpler language that explains a great majority of the examples by using simple rules. This is accomplished by gradually varying the error function in such a way that the net is eventually allowed to classify clearly but incorrectly those strings that are difficult to learn, which are treated as exceptions. The results show that in this way, the DTRNN usually learns a simplified approximate language.

1. Introduction

Discrete-time recurrent neural networks (DTRNN) may be trained to recognize regular languages from sets of example and counterexample strings [2, 10, 4, 12, 6, 3], under the intuitive assumption that, as state machines, DTRNN can emulate deterministic finite automata (DFA), that is, regular-language recognizers. DTRNN can indeed learn simple DFA but they do so at a considerable expense of resources, and their learning behavior does not nicely scale to larger DFA. On the other hand, very efficient discrete algorithms exist [9] that are capable of doing the same task in polynomial time.

One might contend that there is no use in pursuing the study of DTRNN as adaptive recognizers. But if one is interested in obtaining not an *exact* DFA that recognizes all the strings in the learning set, which may be quite complex, but a small, *approximate* DFA which explains all but a few examples, then

the poor scaling behavior of DTRNNs and their bias toward very simple representations, which could be otherwise considered as weaknesses, may be used advantageously to accomplish the task of obtaining a simple and approximate language.

This approach has been explored by Gori and coworkers [5], who view the training set as a noisy or corrupted version of a hypothetical clean learning set for a very simple language; this "noise" only affects the membership of the strings: a few counterexamples are turned into examples and vice versa. When a DTRNN learns to behave as a DFA, its state space tends to organize in clusters of low dimensionality that correspond to the states of the DFA: a clustering algorithm is used in [5] to define these states. A promising alternative approach is described in [7]: an additional term is added to the error function to enforce weight configurations that ensure correct DFA behavior of the DTRNN.

In this work, we propose the use of an error function that varies during the learning process in such a way that the net is eventually allowed to classify clearly but incorrectly those strings that are difficult to learn, which are treated as exceptions.

2. Definitions and methods

Languages and automata: A *language* is a set of finite-length strings over a finite alphabet $\Sigma = \{\sigma_1, \sigma_2, \dots, \sigma_{|\Sigma|}\}$. The length of a string w will be denoted $|w|$. A *deterministic finite automaton* (DFA) is defined as a 5-tuple $M = (Q, \Sigma, \delta, q_I, F)$, where Q is a finite number of *states*, Σ is the input alphabet, δ is the next-state function $\delta : Q \times \Sigma \rightarrow Q$ which defines which state $q' = \delta(q, \sigma)$ is reached by the automaton after reading symbol σ when in state q , $q_I \in Q$ is the initial state of the automaton (before reading any string) and $F \subseteq Q$ is the set of accepting states of the automaton. The language $L(M)$ accepted by the automaton contains all the strings that bring the automaton to an accepting state. The languages accepted by DFA are called *regular languages*.

Second-order DTRNN as string processors: We have used a second-order DTRNN identical to those used in [4] and [3], having n_U inputs and n_X state units. The value of the k -th input at time t is denoted by $u_k[t]$; the state of the i -th state unit *after* processing the input at time t is denoted by $x_i[t]$. The dynamics of the network is defined by $x_i[t] = g\left(\sum_{j=1}^{n_X} \sum_{k=1}^{n_U} W_{ijk} x_j[t-1] u_k[t]\right)$ where $i = 1, \dots, n_X$ and $g(x) = 1/(1 + \exp(-x))$. When a DTRNN is used to process strings, the number of input units n_U is set to $|\Sigma|$, the size of the alphabet and a local (exclusive or *one-hot*) encoding is used so that, for each input symbol, the next state is computed by a different single-layer perceptron. After complete processing of a string w , a value of $x_1[|w|]$ close to 1 denotes acceptance and a value close to 0, rejection; two different criteria are used: the *tight* criterion (rejection: $x_1[|w|] \in]0, 2\tau[$, acceptance: $x_1[|w|] \in [1 - 2\tau, 1[$) and the *loose* criterion (rejection: $x_1[|w|] < 0.5$; acceptance: $x_1[|w|] > 0.5$).

Second-order DTRNN are capable of emulating the behavior of any finite-state machine, provided that $n_X = |Q| + 1$ [8, 1] (although a smaller second-order DTRNN may also emulate the DFA). It may be therefore expected that small DTRNN can only stably emulate and furthermore, can only easily *learn* small DFA, although the nature of the inductive bias of gradient-based learning algorithms is not well known.

Training algorithms and error functions: We use a gradient-based method, the *real-time recurrent learning* (RTRL) [13], but this choice is not critical to our study. We train the weights as in [4] and the initial state $x_i[0]$, $i = 1, \dots, n_X$ as in [3]. In most papers, the error function for each string w is $0.5(x_1[|w|] - T_w)^2$, where the target T_w , is either τ (rejection) or $1 - \tau$ (acceptance). A training run is finished when all $|x_1[|w|] - T_w| < \epsilon$ (here, $\epsilon = \tau$, but any $\epsilon \leq \tau$ is adequate).

The new error function is chosen so that it varies during learning. Consider a learning set in which a large majority of strings are labeled following the rules of a simple DFA whereas a small number of them do not follow these simple rules (their membership labels are reversed). When one presents this learning set to a small second-order DTRNN, it seems that the net initially easily learns the correct behavior for a large majority of strings, but fails to classify clearly and correctly a small number of them; a study of the regions of state space visited by the DTRNN indicates the presence of clusters that may be assumed to correspond to the states of the simple DFA classifying the majority of strings. Later, when trying to learn the labels of this smaller number of strings, the clusters dissolve and the generalization performance (error rate) of the networks when processing longer strings degrades, possibly because the DTRNN is incapable of representing (or learning) a more complex DFA.

Accordingly, we propose the following scheme: the error function will change when the network classifies clearly (even if not correctly) a fraction of the training set whose size is equal (or close) to the expected number of correct strings. The new error function will allow strings to be *incorrectly* but *clearly* classified by having two minima: one at the correct target T_w and another one at the reversed target $1 - T_w$. This is equivalent to performing unsupervised learning at the later stage of the learning process. It is assumed that strings that are already correctly classified will not abandon the minimum around the correct target, whereas strings incorrectly classified will “fall” in either of the two target values, now that they are given an additional choice. The transition to the new function is gradual, to avoid disturbing the gradient descent algorithm. The initial (completely supervised) error function is $e_s(x, T_w) = \frac{1}{2}(x - T_w)^2$ and the final error function (completely unsupervised) is $e_u(x) = \frac{1}{2} \frac{(x-\tau)^2(x-(1-\tau))^2}{(x-\tau)^2+(x-(1-\tau))^2}$ which has two minima at $x = \tau$ and $x = 1 - \tau$ and a maximum at $x = 0.5$. The transition is modelled by using an increasing parameter $\lambda \in [0, 1]$ to combine both functions: $e(x, T_w, \lambda) = (1 - \lambda)e_s(x, T_w) + \lambda e_u(x)$. The transition between one minimum and two minima occurs for $\lambda = \lambda^*(\tau)$; in our experiments, $\lambda^*(0.1) \approx 0.78$. To isolate exceptions, we want this transition to occur

when the net has learned to classify clearly a fraction of the strings similar to the expected fraction of correct strings. For this purpose, if n_{train} is the number of strings in the training set, n_{exc} the expected number of exceptions, and n_{clear} the number of strings that have been clearly classified (n_{clear}), one sets $\lambda = (n_{\text{clear}}/n_{\text{train}})^z$ with z chosen so that $\lambda = \lambda^*$ when $n_{\text{clear}} = n_{\text{train}} - n_{\text{exc}}$.

3. Experiments

Recovering the original language: For each one of Tomita's [11] seven simple languages, a sample of all membership-labeled strings with length 1–9 is generated. Then, it is corrupted (in each simulation) by flipping (reversing) the membership of strings chosen at random according to a uniform distribution with a flip probability r ($r=0.01, 0.02, 0.03$ and 0.05); these strings are marked as exceptions for further analyses. The strings with length 1–6 are used for training and the rest are the test set; n_X is set to the smallest number of hidden units preliminarily observed to converge on sets without exceptions (2 for L_1 ; 3 for L_2, L_3, L_4 and L_7 ; 4 for L_5 and 5 for L_6). For each value of r , 30 runs are performed. The number of exceptions found in each training set (n_{exc}) follows the binomial distribution, which is rather broadly spread around the expected value $\langle n_{\text{exc}} \rangle = rn_{\text{train}}$, with n_{train} the number of strings in the training set.

The learning rate and learning momentum are fixed to $\alpha = 1.0$ and $\eta = 0.2$ respectively, and a run is considered to succeed if the DTRNN recovers the uncorrupted language, that is, it (a) it learns to correctly label (tight criterion) the learning set according to the uncorrupted language; (b) does so in less than 2000 epochs, and (c) correctly labels (loose criterion) the test set. Table 1 shows, for each language and exception rate, the fraction of succeeding jobs, the average number of epochs and its standard deviation. The above criteria for success are very stringent, because they accept only those DTRNN that have captured a stable representation of the original DFA. In a real application, the DFA is not known and exceptions have to be *discovered* by the net by relying on an estimation of the exception rate. Indeed, some jobs isolated a different set of exceptions from the training set (usually less than the ones marked as such), discovering a different DFA with a remarkable generalization performance. Other jobs isolate a very large number of exceptions (25–35), sometimes very early, due to the nature of the stopping criterion (*clear* as opposed to *correct* classification).

Learning an approximate language: The second set of experiments simply probes the ability of the DTRNN to learn simple DFA from the samples. Training is less stringent: a run succeeds if the net learns to label (tight criterion) the training set regardless of the *relabeling* produced, in less than 2000 epochs. To monitor the ability of the DTRNN to learn a simple DFA we use a classical DFA inference algorithm, RPNI, [9], both the corrupted and the DTRNN-re-labeled training set. Table 2 shows, for each language and exception

| | $r = 0.01$ | $r = 0.02$ | $r = 0.03$ | $r = 0.05$ |
|-------|-------------------|-------------------|-------------------|------------------|
| L_1 | 23/30 (26 ± 5) | 17/30 (28 ± 6) | 18/30 (28 ± 6) | 10/30 (36 ± 14) |
| L_2 | 21/30 (180 ± 39) | 17/30 (221 ± 51) | 4/30 (195 ± 21) | 5/30 (250 ± 120) |
| L_3 | 14/30 (300 ± 480) | 9/30 (310 ± 410) | 10/30 (190 ± 360) | 6/30 (250 ± 340) |
| L_4 | 11/30 (160 ± 190) | 3/30 (120 ± 40) | 7/30 (240 ± 300) | 2/30 (80 ± 17) |
| L_5 | 8/30 (570 ± 180) | 6/30 (430 ± 180) | 7/30 (880 ± 670) | 2/30 (190 ± 110) |
| L_6 | 11/30 (140 ± 65) | 10/30 (280 ± 330) | 5/30 (500 ± 720) | 2/30 (148 ± 72) |
| L_7 | 5/30 (159 ± 85) | 9/30 (222 ± 85) | 6/30 (220 ± 100) | 2/30 (381 ± 170) |

Table 1: Ability of the network to recover Tomita's languages from corrupted samples

rate, the number of succeeding jobs and the average and standard deviation of (a) the number of states RPNI produces for the corrupted and the relabeled training set ($|Q|_{\text{before}}$ and $|Q|_{\text{after}}$ respectively); the number n_{exc} of exceptions isolated; and the sum of the two precedent values, $n_{\text{exc}} + |Q|_{\text{after}}$, a measure of approximation (also used in [5]). As can be seen, $|Q|_{\text{after}}$ is usually smaller than $|Q|_{\text{before}}$, and in many cases it approaches the size of the original (Tomita) DFA.

4. Concluding remarks

The results show that DTRNN can be used to infer simple, approximate DFA describing a majority of the examples when the sample cannot be described by a simple DFA due to the presence of noise in the membership labels or to the fact that the sample indeed does not correspond to any simple DFA. Our method, which uses an error function that changes during training to accommodate exceptions, only needs a hint about the expected exception rate in the learning set; when this rate is not known, the parameter may be used to roughly regulate the quality of the approximation. We plan to study the stability of the DFA representation learned by the DTRNN.

References

- [1] R.C. Carrasco, M.L. Forcada, M.A. Valdés, R.P. Neco. Stable encoding of finite-state machines in discrete-time recurrent neural nets with sigmoid units, submitted to *Neural Computation*.
- [2] A. Cleeremans, D. Servan-Schreiber, J.L. McClelland. *Neural Computation*, 1(3):372-381, 1989.
- [3] M.L. Forcada, R.C. Carrasco. *Neural Computation*, 7(5):923-930, 1995.
- [4] C.L. Giles, C.B. Miller, D. Chen, H.H. Chen, G.Z. Sun, Y.C. Lee. *Neural Computation*, 4(3):393-405, 1992.

| L_1 | SUCCESS RATE | $ Q _{\text{before}}$ | $ Q _{\text{after}}$ | n_{exc} | $ Q _{\text{after}+n_{\text{exc}}}$ |
|-------|--------------|-----------------------|----------------------|------------------|-------------------------------------|
| 0.01 | 27/30 | 5.70(2.75) | 1.00(0.00) | 1.11(0.74) | 2.11(0.74) |
| 0.02 | 18/30 | 8.94(2.93) | 1.00(0.00) | 2.00(0.94) | 3.00(0.94) |
| 0.03 | 18/30 | 9.83(2.50) | 1.00(0.00) | 2.83(1.38) | 3.83(1.38) |
| 0.05 | 19/30 | 12.68(1.72) | 1.00(0.00) | 4.63(1.31) | 5.63(1.31) |
| L_2 | SUCCESS RATE | $ Q _{\text{before}}$ | $ Q _{\text{after}}$ | n_{exc} | $ Q _{\text{after}+n_{\text{exc}}}$ |
| 0.01 | 23/30 | 6.35(3.21) | 1.91(0.58) | 1.91(2.22) | 3.83(1.86) |
| 0.02 | 18/30 | 7.50(3.22) | 2.11(0.94) | 2.33(2.08) | 4.44(2.01) |
| 0.03 | 11/30 | 10.45(1.78) | 2.36(1.07) | 3.73(2.05) | 6.09(1.98) |
| 0.05 | 17/30 | 10.94(2.58) | 1.76(0.73) | 4.82(2.31) | 6.59(1.88) |
| L_3 | SUCCESS RATE | $ Q _{\text{before}}$ | $ Q _{\text{after}}$ | n_{exc} | $ Q _{\text{after}+n_{\text{exc}}}$ |
| 0.01 | 14/30 | 6.36(2.38) | 3.00(0.00) | 0.93(0.80) | 3.93(0.80) |
| 0.02 | 14/30 | 9.93(3.43) | 3.00(0.00) | 3.43(7.14) | 6.43(7.14) |
| 0.03 | 16/30 | 11.00(2.87) | 2.94(0.24) | 4.88(9.66) | 7.81(9.42) |
| 0.05 | 8/30 | 15.75(3.73) | 3.00(1.00) | 12.88(15.48) | 15.88(15.24) |
| L_4 | SUCCESS RATE | $ Q _{\text{before}}$ | $ Q _{\text{after}}$ | n_{exc} | $ Q _{\text{after}+n_{\text{exc}}}$ |
| 0.01 | 16/30 | 7.56(3.24) | 2.94(1.03) | 5.00(10.42) | 7.94(9.70) |
| 0.02 | 14/30 | 9.29(3.79) | 2.86(0.52) | 4.14(8.64) | 7.00(8.13) |
| 0.03 | 15/30 | 13.40(3.18) | 2.67(1.25) | 11.87(13.27) | 14.53(12.28) |
| 0.05 | 8/30 | 11.88(2.76) | 2.50(0.87) | 12.00(13.64) | 14.50(12.78) |
| L_5 | SUCCESS RATE | $ Q _{\text{before}}$ | $ Q _{\text{after}}$ | n_{exc} | $ Q _{\text{after}+n_{\text{exc}}}$ |
| 0.01 | 7/30 | 5.29(1.48) | 4.00(0.00) | 0.43(0.49) | 4.43(0.49) |
| 0.02 | 6/30 | 7.83(3.48) | 5.00(1.41) | 0.83(0.90) | 5.83(1.77) |
| 0.03 | 6/30 | 12.33(2.62) | 5.50(1.50) | 2.33(0.75) | 7.83(1.95) |
| 0.05 | 1/30 | 17.00(0.00) | 2.00(0.00) | 22.00(0.00) | 24.00(0.00) |
| L_6 | SUCCESS RATE | $ Q _{\text{before}}$ | $ Q _{\text{after}}$ | n_{exc} | $ Q _{\text{after}+n_{\text{exc}}}$ |
| 0.01 | 20/30 | 7.05(2.96) | 4.20(2.09) | 3.85(10.07) | 8.05(9.95) |
| 0.02 | 14/30 | 9.07(2.79) | 5.64(3.35) | 0.93(0.80) | 6.57(3.06) |
| 0.03 | 12/30 | 11.83(2.44) | 7.00(4.34) | 3.50(4.92) | 10.50(5.91) |
| 0.05 | 13/30 | 15.54(3.39) | 6.15(3.80) | 7.23(8.55) | 13.38(8.00) |
| L_7 | SUCCESS RATE | $ Q _{\text{before}}$ | $ Q _{\text{after}}$ | n_{exc} | $ Q _{\text{after}+n_{\text{exc}}}$ |
| 0.01 | 14/30 | 7.36(2.64) | 3.36(1.23) | 7.00(12.20) | 10.36(10.97) |
| 0.02 | 16/30 | 11.12(2.74) | 3.69(1.45) | 6.62(10.93) | 10.31(9.65) |
| 0.03 | 14/30 | 13.29(3.06) | 3.29(1.83) | 12.79(13.93) | 16.07(12.24) |
| 0.05 | 7/30 | 15.86(3.23) | 2.71(1.48) | 18.57(14.07) | 21.29(12.84) |

Table 2: Neural relabeling: success rate, size of exact and approximate DFA, and number of exceptions (see text)

- [5] M. Gori, M. Maggini, G. Soda. *IEEE Transactions on Neural Networks* 9(3):571-575, 1998.
- [6] P. Manolios, R. Fanelli. *Neural Computation* 6(6):1154-1172, 1994.
- [7] R.P. Ñeco, S.C. Kremer, M.L. Forcada. In *Proc. ICANN'98*, v. 2, p. 529-534.
- [8] C.W. Omlin and C.L. Giles. *Journal of the ACM* 43(6):937-972, 1966.
- [9] J. Oncina, P. Garcia. In *Advances in Structural and Syntactic Pattern Recognition*. World Scientific, 1992.
- [10] J.B. Pollack. *Machine Learning* 7:227-252, 1991.
- [11] M. Tomita. In *Proc. 4th Ann. Cognitive Sci. Conf*, p. 105-108.
- [12] R.L. Watrous, G.M. Kuhn, *Neural Computation* 4(3):406-414, 1992.
- [13] R.J. Williams, D. Zipser, *Neural Computation* 1(2):270-280, 1989.