

Can Neural Networks Learn Symbolic Rewriting?

Abstract

This work investigates if the current neural architectures are adequate for learning symbolic rewriting. Two kinds of data sets are proposed for this research -- one based on automated proofs and the other being a synthetic set of polynomial terms. The experiments with use of the current neural machine translation models are performed and its results are discussed. Ideas for extending this line of research are proposed and its relevance is motivated.

Introduction

Neural networks (NNs) turned out to be very useful in several domains. In particular, one of the most spectacular advances achieved with use of NNs has been natural language processing. One of the tasks in this domain is translation between natural languages – neural machine translation (NMT) systems established here the state-of-the-art performance. Recently, NMT produced first encouraging results in the autoformalization task BIBREF0, BIBREF1, BIBREF2, BIBREF3 where given an informal mathematical text in the goal is to translate it to its formal (computer understandable) counterpart. In particular, the NMT performance on a large synthetic -to-Mizar dataset produced by a relatively sophisticated toolchain developed for several decades BIBREF4 is surprisingly good BIBREF3, indicating that neural networks can learn quite complicated algorithms for symbolic data. This inspired us to pose a question: Can NMT models be used in the formal-to-formal setting? In particular: Can NMT models learn symbolic rewriting?

The answer is relevant to various tasks in automated reasoning. For example, neural models could compete with symbolic methods such as inductive logic programming BIBREF5 (ILP) that have been

previously experimented with to learn simple rewrite tasks and theorem-proving heuristics from large formal corpora BIBREF6. Unlike (early) ILP, neural methods can however easily cope with large and rich datasets, without combinatorial explosion.

Our work is also an inquiry into the capabilities of NNs as such, in the spirit of works like BIBREF7.

Data

To perform experiments answering our question we prepared two data sets – the first consists of examples extracted from proofs found by ATP (automated theorem prover) in a mathematical domain (AIM loops), whereas the second is a synthetic set of polynomial terms.

Data :: The AIM data set

The data consists of sets of ground and nonground rewrites that came from Prover9 proofs of theorems about AIM loops produced by Veroff BIBREF8.

Many of the inferences in the proofs are paramodulations from an equation and have the form $s = t$

$u[(s)] = v[u[(t)]] = v$ where s, t, u, v are terms and θ is a substitution. For the most common equations $s = t$, we gathered corresponding pairs of terms $(u[\theta(s)], u[\theta(t)])$ which were rewritten from one to another with $s = t$. We put the pairs to separate data sets (depending on the corresponding $s = t$): in total 8 data sets for ground rewrites (where θ is trivial) and 12 for nonground ones. The goal will be to learn rewriting for each of this 20 rules separately.

Terms in the examples are treated as linear sequences of tokens where tokens are single symbols

(variable / constant / predicate names, brackets, commas). Numbers of examples in each of the data sets vary between 251 and 34101. Lengths of the sequences of tokens vary between 1 and 343, with mean around 35. These 20 data sets were split into training, validation and test sets for our experiments (60 %, 10 %, 30 %, respectively).

In Table TABREF4 and Table TABREF5 there are presented examples of pairs of AIM terms in TPTP BIBREF9 format, before and after rewriting with, respectively, ground and nonground rewrite rules.

Data :: The polynomial data set

This is a synthetically created data set where the examples are pairs of equivalent polynomial terms. The first element of each pair is a polynomial in an arbitrary form and the second element is the same polynomial in a normalized form. The arbitrary polynomials are created randomly in a recursive manner from a set of available (non-nullary) function symbols, variables and constants. First, one of the symbols is randomly chosen. If it is a constant or a variable it is returned and the process terminates. If a function symbol is chosen, its subterm(s) are constructed recursively in a similar way.

The parameters of this process are set in such a way that it creates polynomial terms of average length around 25 symbols. Terms longer than 50 are filtered out. Several data sets of various difficulty were created by varying the number of available symbols. This was quite limited – at most 5 different variables and constants being a few first natural numbers. The reason for this limited complexity of the input terms is because normalizing even a relatively simple polynomial can result in a very long term with very large constants – which is related especially to the operation of exponentiation in polynomials.

Each data set consists of different 300 000 examples, see Table TABREF7 for examples. These data sets were split into training, validation and test sets for our experiments (60 %, 10 %, 30 %, respectively),

respectively).

Experiments

For experiments with both data sets we used an established NMT architecture BIBREF10 based on LSTMs (long short-term memory cells) and implementing the attention mechanism.

After a small grid search we decided to inherit most of the hyperparameters of the model from the best results achieved in BIBREF3 where -to-Mizar translation is learned. We used relatively small LSTM cells consisting of 2 layers with 128 units. The “scaled Luong” version of the attention mechanism was used, as well as dropout with rate equal 0.2. The number of training steps was 10000. (This setting was used for all our experiments described below.)

Experiments :: AIM data set

First, NMT models were trained for each of the 20 rewrite rules in the AIM data set. It turned out that the models, as long as the number of examples was greater than 1000, were able to learn the rewriting task very well, reaching 90% of accuracy on separated test sets. This means that the task of applying single rewrite step seems relatively easy to learn by NMT. See Table TABREF11 for all the results.

We also run an experiment on the joint set of all rewrite rules (consisting of 41396 examples). Here the task was more difficult as a model needed not only to apply rewriting correctly, but also choose “the right” rewrite rule applicable for a given term. Nevertheless, the performance was also very good, reaching 83% of accuracy.

Experiments :: Polynomial data set

Then experiments on more challenging but also much larger data sets for polynomial normalization were performed. Depending on the difficulty of the data, accuracy on the test sets achieved in our experiments varied between 70% and 99%. The results in terms of accuracy are shown in Table TABREF13.

This high performance of the model encouraged a closer inspection of the results. First, we checked if in the test sets there are input examples which differs from these in training sets only by renaming of variables. Indeed, for each of the data sets in test sets are 5 - 15 % of such “renamed” examples. After filtering them out the measured accuracy drops – but only by 1 - 2 %.

An examination of the examples wrongly rewritten by the model was done. It turns out that the wrong outputs almost always parse (in 97 - 99 % of cases they are legal polynomial terms). Notably, depending on the difficulty of the data set, as much as 18 - 64 % of incorrect outputs are wrong only with respect to the constants in the terms. (Typically, NMT model proposes too low constants compared to the correct ones.) Below 1 % of wrong outputs are correct modulo variable renaming.

Conclusions and future work

NMT is not typically applied to symbolic problems, but surprisingly, it performed very well for both described tasks. The first one was easier in terms of complexity of the rewriting (only one application of a rewrite rule was performed) but the number of examples was quite limited. The second task involved more difficult rewriting – multiple different rewrite steps were performed to construct the examples. Nevertheless, provided many examples, NMT could learn normalizing polynomials.

We hope this work provides a baseline and inspiration for continuing this line of research. We see several interesting directions this work can be extended.

Firstly, more interesting and difficult rewriting problems need to be provided for better delineation of the strength of the neural models. The described data are relatively simple and with no direct relevance to the real unsolved symbolic problems. But the results on these simple problems are encouraging enough to try with more challenging ones, related to real difficulties – e.g. these from TPDB data base.

Secondly, we are going to develop and test new kinds of neural models tailored for the problem of comprehending symbolic expressions. Specifically, we are going to implement an approach based on the idea of TreeNN, which may be another effective approach for this kind of tasks BIBREF7, BIBREF12, BIBREF13. TreeNNs are built recursively from modules, where the modules corresponds to parts of symbolic expression (symbols) and the shape of the network reflects the parse tree of the processed expression. This way model is explicitly informed on the exact structure of the expression, which in case of formal logic is always unambiguous and easy to extract. Perhaps this way the model could learn more efficiently from examples (and achieve higher results even on the small AIM data sets). The authors have a positive experience of applying TreeNNs to learn remainders of arithmetical expressions modulo small natural numbers – TreeNNs outperformed here neural models based on LSTM cells, giving almost perfect accuracy. However, this is unclear how to translate this TreeNN methodology to the tasks with the structured output, like the symbolic rewriting task.

Thirdly, there is an idea of integrating neural rewriting architectures into the larger systems for automated reasoning. This can be motivated by the interesting contrast between some simpler ILP systems suffering for combinatorial explosion in presence of a large number of examples and neural methods which definitely benefit from large data sets.

We hope that this work will inspire and trigger a discussion on the above (and other) ideas.

Acknowledgements

Piotrowski was supported by the grant of National Science Center, Poland, no. 2018/29/N/ST6/02903, and by the European Agency COST action CA15123. Urban and Brown were supported by the ERC Consolidator grant no. 649043 AI4REASON and by the Czech project AI&Reasoning CZ.02.1.01/0.0/0.0/15_003/0000466 and the European Regional Development Fund. Kaliszyk was supported by ERC Starting grant no. 714034 SMART.