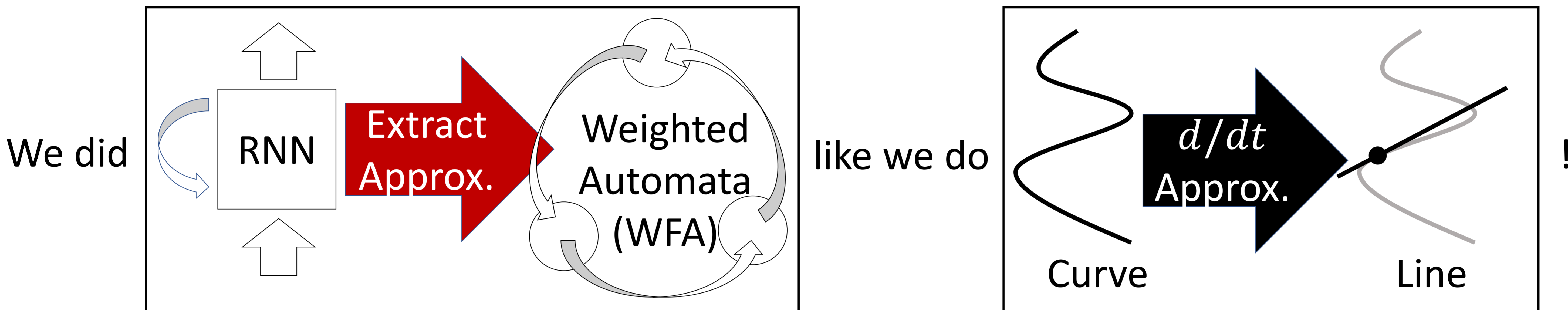


Weighted Automata Extraction from Recurrent Neural Networks via Regression on State Spaces



Takama Okudono, Masaki Waga, Taro Sekiyama, Ichiro Hasuo
@ National Institute of Informatics & the Graduate University for Advanced Studies (SOKENDAI), Japan

What's this?



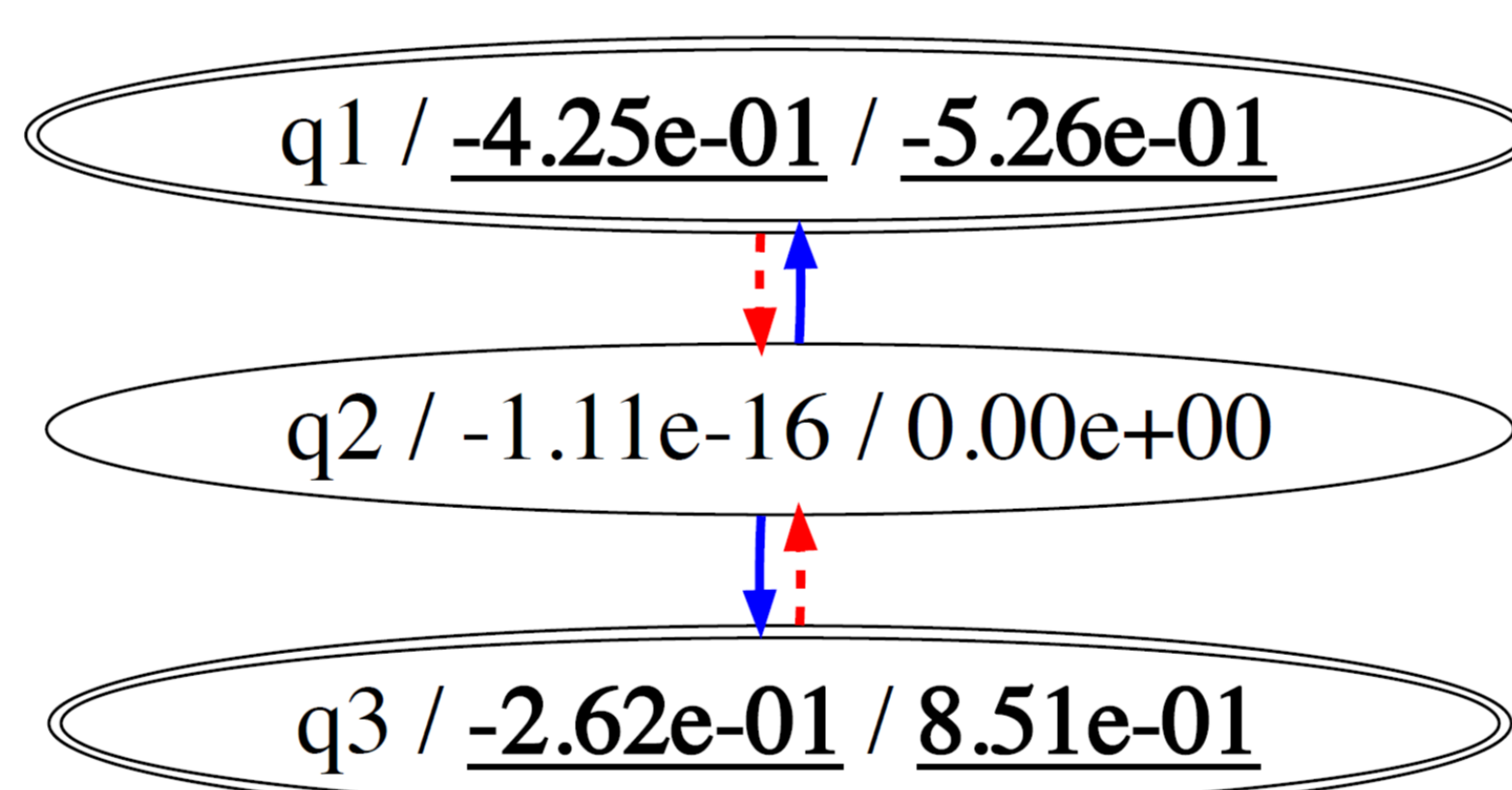
Experimental results

(1) Extraction is quick and accurate!

(2) Extraction yields interpretability!

(3) Extraction makes the inference faster!

(Σ , Q_A)	RGR(2)	RGR(5)	BFS(500)	BFS(1000)	BFS(2000)	BFS(3000)	BFS(5000)
(4, 10)	2.17 / 286	2.39 / 338	26.8 / 165	9.77 / 279	4.36 / 545	4.07 / 716	2.33 / 1390
(6, 10)	2.45 / 1787	2.54 / 1302	6.99 / 386	4.48 / 641	4.08 / 1218	3.15 / 1410	2.28 / 2380
(10, 10)	4.68 / 7462	4.46 / 5311	22.5 / 928	11.9 / 1562	5.90 / 3521	4.55 / 3638	3.55 / 5571
(10, 15)	5.62 / 8941	5.78 / 8564	21.2 / 2155	10.6 / 4750	7.87 / 5692	5.71 / 7344	5.27 / 7612
(10, 20)	3.70 / 7610	3.79 / 7799	6.24 / 2465	10.1 / 2188	6.13 / 3106	3.70 / 5729	3.63 / 7473
(15, 10)	7.34 / 9569	5.52 / 10000	13.5 / 3227	8.01 / 6765	6.07 / 7916	5.98 / 8911	6.17 / 8979
(15, 15)	8.44 / 10000	5.58 / 9981	16.3 / 2675	9.24 / 4850	7.28 / 5135	9.88 / 7204	6.44 / 8425
(15, 20)	9.16 / 7344	5.15 / 7857	13.7 / 2224	7.26 / 3823	6.60 / 5744	4.96 / 5674	4.01 / 9464
Total	5.45 / 6625	4.40 / 6394	15.9 / 1778	8.92 / 3107	6.04 / 4110	5.25 / 5078	4.21 / 6549



x1300 faster
by using WFA
Instead of RNN

- We trained RNNs from randomly generated WFAs and applied the extraction to the RNNs. The cells show “{error}/time.”
- Rows show various settings of the random WFAs / Columns show extraction methods: RGR(*) are ours, and BFS(*) are naïve baselines.
- Two tables use different ways to generate WFAs (see our paper!)

- We trained RNNs from a weighted language (see “Math stuff”) that are so complex that any WFAs cannot express it and applied our extraction.
- The extracted WFAs represent a simplified version of the original language, and it can be graphically representable as above.

- We used the same RNNs and extracted WFAs and compared inference time for random inputs.

Method

Algorithm 1 Answering equivalence queries

```

1: procedure ANS-EQ
   Input: RNN  $R = (\alpha_R, \beta_R, g_R)$ , WFA  $A = (Q_A, \alpha_A, \beta_A, (A_\sigma)_{\sigma \in \Sigma})$ , error tolerance  $e > 0$  and concentration threshold  $M \in \mathbb{N}$ 
   Output: a counterexample, or Equivalent
2: Initialize  $p: \mathbb{R}^d \rightarrow \mathbb{R}^{Q_A}$  so that  $p(\delta_R(\varepsilon)) = \delta_A(\varepsilon)$ 
3: queue  $\leftarrow \{\varepsilon\}$ ; visited  $\leftarrow \emptyset$ 
4: while queue is non-empty do
5:    $h \leftarrow \text{pop}(\text{queue})$ 
   ▷ Pop the element of the maximum priority
6:   if  $|f_R(h) - f_A(h)| \geq e$  then
7:     return  $h$  ▷ return a counterexample
8:   result  $\leftarrow \text{CONSISTENT?}(h, \text{visited}, p)$ 
9:   if result = NG then
10:    learn  $p$  by regression, so that  $p(\delta_R(h')) = \delta_A(h')$  holds for all  $h' \in \text{visited} \cup \{h\}$ 
11:    visited  $\leftarrow \text{visited} \cup \{h\}$ 
12:    visited'  $\leftarrow p(\delta_R(\text{visited}))$ 
13:    #  $vn \leftarrow \{x \in \text{visited}' \mid x \simeq_A p(\delta_R(h))\}$ 
14:    if #  $vn \leq M$  then
15:      pr  $\leftarrow \min_{h' \in \text{visited}' \setminus \{h\}} d(p(\delta_R(h)), p(\delta_R(h')))$ 
      ▷  $d$  is the Euclidean distance
16:      push  $h\sigma$  to queue with priority pr for  $\sigma \in \Sigma$ 
17: return Equivalent

```

- Balle and Mohri's WFA learning algorithm reduces the problem of checking whether the target RNN and a given WFA are equivalent in respect of behavior.
- The algorithm is on the left. The main ideas are:
 - Run best-first search on Σ^* and find the difference between the RNN and the WFA (Line 4-7, 16)
 - Find the relation between the RNN and the WFA by making a map from the internal state of RNN to the internal state of WFA by GPR (Line 10)
 - Prioritize the visit where “the search is not enough.” (Line 5, 14-15)
 - Prune the searching where “the search is enough.” (Line 12-13)

Math stuff...

[Weighted language] Function $\Sigma^* \rightarrow \mathbb{R}$
[WFA] WFA A is a quadruple of

- Q_A — states
- $\alpha_A \in \mathbb{R}^{Q_A}$ — initial vector
- $\beta_A \in \mathbb{R}^{Q_A}$ — final vector
- $(A_\sigma)_{\sigma \in \Sigma} \in \mathbb{R}^{Q_A \times Q_A}$ — transition matrix

WFA induces its configuration δ_A and output f_A by:

- $\delta_A(w_1 \dots w_n) = \alpha_A^T A_{w_1} \dots A_{w_n}$
- $f_A(w_1 \dots w_n) = \delta_A(w_1 \dots w_n) \beta_A$

[RNN] RNN R is a triple of

- $\alpha_R \in \mathbb{R}^d$ — initial state
- $\beta_R \in \mathbb{R}^d \rightarrow \mathbb{R}$ — final state
- $g_R: \mathbb{R}^d \times \Sigma^* \rightarrow \mathbb{R}^d$ — transition func.

RNN induces its configuration δ_R and output f_A by:

- $\delta_R(w_1 \dots w_n) = \begin{cases} g_R(\delta_R(w_1 \dots w_{n-1}), w_n) & n > 0 \\ \alpha_R & n = 0 \end{cases}$
- $f_R(w_1 \dots w_n) = \beta_R(\delta_R(w_1 \dots w_n))$

Future work

- Improving scalability
Currently, $|\Sigma| = 15$ is the maximum. Can we make it work for larger $|\Sigma|$ so that it is applicable for NLP?
- Giving a theoretical guarantee
It is a heuristics now, and there is no guarantee for the superiority of our method nor termination.
- Applying for quality assurance of RNN
Model-checking is a field to generate the proof of the (typically) safety of the system. Can we combine our technique with model checking of WFAs and make a technique for the quality assurance of RNN?
- Forcing the WFA to be probabilistic
Even if the target system outputs the values in $[0, 1]$, Balle and Mohri's algorithm does not necessarily output WFAs whose output is in it. Can we fix it and improve the applicability to the RNNs whose outputs are probabilistic?