Finite State Transducers for Information Extraction TID 2009/10

Marek Schmidt

FIT BUT

2009-12-31

Marek Schmidt (FIT BUT) Finite State Transducers for Information Extract

2009-12-31 1 / 14

Introduction

We want to show that a practical Information Extraction system can be implemented as a Finite State Transducer. For that we need to...

- Define a string representation of text with annotations.
- Operators on FSTs.
- Implement Information Extraction system using this set of operators.

Finite State Transducer

$$T = (Q, \Sigma, \Gamma, I, F, \delta)$$
(1)

- Q ... Finite set of states. (2)
- Σ ... Finite set, input alphabet. (3)
- Γ ... Finite set, output alphabet(.4)
- $I \subseteq Q$... Initial states (5)
- $F \subseteq Q$... Final states (6)
- $\delta \subseteq \boldsymbol{Q} \times (\boldsymbol{\Sigma} \cup \{\epsilon\}) \times (\boldsymbol{\Gamma} \cup \{\epsilon\}) \times \boldsymbol{Q} \quad \dots \quad \text{Transitions}$ (7)

$$x[T]y \subseteq \Sigma^* \times \Gamma^* \dots$$
 Behavior

(8)

Operations on Finite State Transducers

Atom	<i>x</i> / <i>y</i>	x[x/y]y	(9)
Identity	R	$x[R]x$ iff $x \in R$	(10)
Union	$T \cup U$	$x[T \cup U]y$ iff $x[T]y$ or $x[U]y$	(11)
Intersection	$T \cap U$	$x[T \cap U]y$ iff $x[T]y$ and $x[U]y$	(12)
Concatenation	T.U	xw[T.U]yz iff $x[T]y$ and $w[U]z$	(13)
Iteration	T *	Minimal such as $\epsilon[\mathcal{T}^*]\epsilon$ and	(14)
		$x[T^*]y \wedge w[T]z \Rightarrow xw[T^*]yz.$	
Composition	$T \circ U$	$x[T \circ U]y$ iff $\exists z.x[T]z$ and $z[U]y$	(15)

- 3 >

Composition (*e*-free)

$$T = (Q_T, \Sigma_T, \Gamma_T, I_T, F_T, \delta_T)$$
(16)

$$U = (Q_U, \Sigma_U, \Gamma_U, I_U, F_U, \delta_U)$$
(17)

$$T \circ U = (Q_V, \Sigma_V, \Gamma_V, I_V, F_V, \delta_V)$$
(18)

•
$$Q_V = Q_T \times Q_U$$

- $\Sigma_V = \Sigma_T$
- $\Gamma_V = \Gamma_U$

•
$$I_V = \{(t, u) \in Q_V | t \in I_T \land u \in I_U\}$$

• $F_V = \{(t, u) \in Q_V | t \in F_T \land u \in F_U\}$

$$\delta_{V} = \{ ((q_{T}, q_{U}), s_{T}, t_{V}, (r_{T}, r_{U})) \in Q_{V} \times \Sigma_{V} \times \Gamma_{V} \times Q_{V} | \quad (19) \\ \exists x. (q_{T}, s_{T}, x, r_{T}) \in \delta_{T} \land (q_{U}, x, t_{V}, r_{U}) \in \delta_{V}) \}$$

イロン イ理 とく ヨン ト ヨン・

Information Extraction

- Extract structured information from unstructured natural language text.
- Entity Recognition
- Relationship Extraction

by the means of patterns on annotated text.

	Oracle	has	acquired	Sun	Microsystems	for	7.4	billion.
POS	Ν	V	V	Ν	N	Ι	С	С
PHRASE	N		V		Ν	Ι		С
ROLE	Recipient				Theme		Мо	ney
EVENT				Acqu	isition			

Annotation Model

L... Finite set of annotation labels

- Coloured strings
 - Extend alphabet $\Sigma' = \Sigma \times 2^{L}$
 - FST with Predicates and Identitites (Noord01)

2 Markup

• Extend alphabet $\Sigma' = \Sigma \cup (L \times \{_{begin, end}\})$

The Markup route seems more practical...

Operations on labels

$$\begin{array}{ll} \text{Match} & \langle x \rangle & \equiv x_{begin}.(\Sigma - \{x_{begin}, x_{end}\})^*.x_{end} & (20) \\ \text{Label} & T:x & \equiv \epsilon/x_{begin}.T.\epsilon/x_{end} & (21) \\ \text{Context} & U = T_M & \delta_U = \delta_T \cup \{(q, l, l, q)| & (22) \\ & q \in Q_U, l \in (L-M) \times \{_{begin}, end\} \} \end{array}$$

(see Ignore operator in (KaplanAndKay94))
 Example

$$(\langle D \rangle.(\langle A \rangle)^*.(\langle N \rangle)^* \langle N \rangle) : NP$$
(23)

 $((\langle NP \rangle) : Recipient.(has acquired)_{\emptyset}.(\langle NP \rangle) : Theme) : Acquisition$ (24)

Information Extraction Pipeline

Composition of transducers.

- Sentence splitting
 - Dummy label to ends of sentences to act as a sentinel.
- 2 Tokenization
 - w_{begin} lorem w_{end} w_{begin} ipsum w_{end}
- Morphology analysis
 - Lemmatization
 - Stemming
- Part-of-speech tagging
- Gazetteers
- Shallow Parsing
- Extraction Patterns

Morphology example



do/do∪does/do∪done/do∪did/do∪ see/see∪saw/see∪seen/see∪ saw/saw∪sawn/saw∪ dot/dot∪dotted/dot

Finite State Transducers for Information Extrac

Efficiency discussion

- Unlike FSM, no equivalent DFST for every NFST.
 - But we need nondeterminism in NLP anyway for ambiguities.
- Subsequential transducers
 - FST with deterministic input.
- Finitely subsequential transducers
 - deterministic transducers with finite set of final output strings.
- Finding equivalent finitely subsequential transducer is decidable (Mohri2003)

2009-12-31

12/14

- *Label* operation in general leads to non-subsequentiable transducers.
- on-line implementation of composition operation.

Conclusion

- defined a set of operators on annotated text.
- outlined an implementation of all necessary components of an information extraction system.
- implemented a proof-of-concept on top FSA Utilities toolbox http://www.let.rug.nl/~vannoord/Fsa/ (Noord99)
- still more theoretical than practical.

References

- (KaplanAndKay94) R. Kaplan and M. Kay. 1994. Regular Models of Phonological Rule System, *Computational Linguistics*
- (Mohri2003) C. Allauzen and M. Mohri. 2003. Finitely Subsequential Transducers, *International Journal of Foundations* of Computer Science
- (Noord99) G. van Noord and D. Gerdemann. 1999. An Extendible Regular Expression Compiler for Finite-state Approaches in Natural Language Processing. WIA 99
- (Noord01) G. van Noord and D. Gerdemann. 2001. Finite State Transducers with Predicates and Identities. *Grammars*

→ ∃ →