















The regular languages A ⊆ Σ\* form a Boolean Algebra



• Since they are closed under intersection and complement.

6

Informatics 1 School of Informatics, University of Edinburgh

### Determinism



Can always convert to an equivalent DFA for which every state has exactly one transition leaving the state for each input symbol.

For this machine there is exactly one trace for each input string

Add a new "black hole" state, •

Proof

- For every pair (s, a) for which there is no state t with a transition T(s, a, t), add a transition T(s, a, •). This includes a transition T(•, a, •) for each  $a \in \Sigma$ . You cannot escape from the black hole.

7

- The black hole is not an accepting state.
- This machine accepts the same language as the original.

Informatics 1 School of Informatics, University of Edinburgh

Non Determ	inis	sm	.0.7	N N N N N N N N N N N N N N N N N N N
In a non-deterministic machine (NFA any number of transitions with the s leaving to different successor states	A), eacl ame in	h state put sy	e may mbol,	have
$\int_{0}^{0} \frac{1}{1} \frac{0}{1}$		•		ſ
$\longrightarrow 0$ (1) (2)		0	1	
1 1	0	0	0,1	
	1	2		
	2			
				1
Informatics 1 School of Informatics, University of Edinburgh				











Internal Transit	ic	n	S	N N N N N N N N N N N N N N N N N N N	- 4 · H
We sometimes add <b>internal transitions</b> – labelled $\varepsilon$ – to a non-deterministic machine (NFA). This is a state change that consumes po input	0	<b>0</b> 0 2	<mark>1ε</mark> 10		
It introduces non-determinism in the observed behaviour of the machine.	2		0		
$ \xrightarrow{ 0 \ 1 \ 0 \ 1 \ 0 \ 0 \ \epsilon } $	0 1 2	0ε* 0 2	<b>1ε*</b> 1,0		
Informatics 1 School of Informatics, University of Edinburgh 12					

































If R and S are regular expressions then the equation  $X = R \mid X S$ has a solution X = R S\* If  $\varepsilon \notin L(S)$  then this solution is unique.  $L_2 = b \mid L_2 (b b b \mid b a)$  $L_2 = b (b b b \mid b a)^*$ 

# regular expressions

- any character is a regexp
- matches itself • if R and S are regexps, so is RS
- matches
  a match for R followed by a match for S • if R and S are regexps, so is RIS
- matches any match for R or S (or both)
- if R is a regexp, so is R\*
- matches
  any sequence of 0 or more matches for R The algebra of regular expressions also includes elements 0 and 1
  - 0 matches nothing; 1 matches the empty string







### **REs** and **FSAs**

- Regular expressions can be viewed as a textual way of specifying the structure of finite-state automata
- Finite-state automata are a way of implementing regular expressions
- Regular expressions denote regular sets of strings each regular set is recognised by some FSA

# Regular expressions

- A formal language for specifying text strings
- How can we search for any of these?
  - woodchuck
  - woodchucks
  - Woodchuck
  - Woodchucks



### Regular Expressions for Textual Searches

#### Who does it?

Everybody:

- Web search engines, CGI scripts
- Information retrieval
- Word processing (Emacs, vi, MSWord)
- Linux tools (sed, awk, grep)
- Computation of frequencies from corpora
- Perl





# **Regular Expression**

- **Regular expression:** formula in algebraic notation for specifying a set of strings
- String: any sequence of alphanumeric characters -letters, numbers, spaces, tabs, punctuation marks
- Regular expression search
  - -pattern: specifying the set of strings we want to search for

-corpus: the texts we want to search through

#### Basic Regular Expression Patterns

- Case sensitive: d is not the same as D
- Disjunctions: [dD] [0123456789]
- Ranges: [0-9] [A-Z]
- Negations: [^Ss] (only when ^ occurs immediately after [)
- Optional characters: ? and \*
- Wild : .
- Anchors: ^ and \$, also \b and \B
- Disjunction, grouping, and precedence: | (pipe)

#### Caret for negation, ^, or anchor

RE	Match (single characters)	Example Patterns Matched
[^A-Z]	not an uppercase letter	"Oyfn pripetchik"
[^Ss]	neither 'S' nor 's'	"I have no exquisite reason for 't"
[^\.]	not a period	"our resident Djinn"
[e/]	either 'e' or '^'	"look up _ now"
a^b	the pattern 'a^b'	"look up <u>a^b</u> now"
^T	T at the beginning of a line	"The Dow Jones closed up one"

### Optionality and Counters

RE	Match	Example Patterns Matched
woodchucks?	woodchuck or woodchucks	"The woodchuck hid"
colou?r	color or colour	"comes in three colours"
(he) {3}	exactly 3 "he"s	"and he said hehehe."

? zero or one occurrences of previous char or expression

- \* zero or more occurrences of previous char or expression
- + one or more occurrences of previous char or expression
- {n} exactly n occurrences of previous char or expression
- $\{n,\,m\}$  between n to m occurrences
- {n, } at least n occurrences

#### Wild card '.'

_		
beg.n a	my char between beg and n	begin, beg'n, begun
big.*dog 1	find lines where big and	the big dog bit the little
(	dog occur	the big black dog bit the

* ? \$ [] [^] [a-z] \ \ \	any character (but newline) previous character or group, repeated 0 or more time previous character or group, repeated 1 or more time previous character or group, repeated 0 or 1 time start of line end of line any character between brackets any character not in the brackets any character between a and z prevents interpretation of following special char or word constituent
\b	word boundary
\{3\} \{3,\} \{3,6	previous character or group, repeated 3 times previous character or group, repeated 3 or more times } previous character or group, repeated 3 to 6 times