

Regular Expressions Example: You can follow this "naturally" without consciously applying the formal definitions. <sup>①</sup>

I: How to represent <sup>Natural</sup> Whole Numbers? IN

Alphabet  $DIG = \{0, 1, \dots, 9\}$

Idea 1:  $DIG^*$  As in lecture, the  $*$  means "zero or more"  
Problem: allows the empty string  $\epsilon$ . "Invisible zero".

Idea 2:  $DIG^+$  : Superscript  $+$  means "one or more".

Second Issue: Allows redundant codes for integers with leading 0s.  
eg.  $007 = 7 = 0000000007$ .

Idea 3: Require the first digit to be 1 thru 9.

$(1 \cup 2 \cup 3 \cup 4 \cup 5 \cup 6 \cup 7 \cup 8 \cup 9) \cdot DIG^*$   
↑  
or UNIX: shorthand  $[1-9]$  followed by

This simple example already uses all three regular operations  $\cup, \cdot, *$ .

Problem: 0 has disappeared.

Idea 4: Union it as a special case:

Naive  $= [1-9] \cdot DIG^* \cup 0$ .

I Now since we allow negative numbers, we want to use an operator  $(-)$  — sign.

"Notation for Optional": "BNF Grammar" [ ] (2)  
 UNIX/Linux -?  
 Regexp/Text ("-" ∪ ε)

("-" ∪ ε) ([1-9] DIG<sup>0</sup> ∪ 0) allows "-0"

(- ∪ ε) [1-9] DIG<sup>0</sup> ∪ 0 does not allow "-0"

<sup>vi</sup> (- ∪ + ∪ ε) [1-9] DIG<sup>0</sup> ∪ 0 allows optional + sign too.

III: Now how about Floating point numbers?

INT • DIG<sup>+</sup>

now leading 0s are needed in the floating part, but trailing 0s cause redundancy, OK.

Issue: can't write .5  
 must write 0.5.

Bug or feature?

If you write (- ∪ + ∪ ε) DIG<sup>0</sup> . DIG<sup>0</sup> you can do .5,

but this also gets "5." or just "."  
 OK cause we can't allow the latter, so let's do

INT . DIG<sup>0</sup> ∪ (- ∪ + ∪ ε) DIG<sup>0</sup> . DIG<sup>+</sup>

Formal defn & concatenation:

(3)

$$A \circ B = \{x \cdot y : x \in A \wedge y \in B\}$$

Contrast with Cartesian Product:

$$A \times B = \{(x, y) : x \in A \wedge y \in B\}$$

$$A = \{a, ab\} \quad A \cdot B = \{a \cdot b, a \cdot bb, ab \cdot b, ab \cdot bb\}$$
$$B = \{b, bb\} \quad = \{ab, abb, abbb\} \text{ size } 3$$

$$A \times B = \{(a, b), (a, bb), (ab, b), (ab, bb)\}$$

always  $|A \times B| = |A| \cdot |B|$ . here, = 4

$$A \circ A = \{x \cdot y : x \in A \wedge y \in A\}$$

$$\neq \{x \cdot x : x \in A\}$$

which would be  $\subseteq$  the  
Doubled language (not lecture)

$$A^1 = A, \quad A^2 = A \circ A, \quad A^3 = A \circ A \circ A \dots \text{ (happily, } \circ \text{ is associative)}$$

Rule of Exponents:  $A^q \circ A^r = A^{q+r}$  like  $0^0 = 1$ .

Needs  $A^0 \equiv \{\epsilon\}$  not  $\emptyset$ . Even  $\emptyset^0 = \{\epsilon\}$  not  $\emptyset$ !