





Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 200

Real Literals in Pascal • Notation: - l stands for $\mathbf{a} + \mathbf{b} + \dots + \mathbf{z} + \mathbf{A} + \mathbf{B} + \dots + \mathbf{Z}$ - d stands for $\mathbf{0} + \mathbf{1} + \mathbf{2} + \dots + \mathbf{9}$ - s stands for $\mathbf{0} + \mathbf{1} + \mathbf{2} + \dots + \mathbf{9}$ - s stands for "sign" (shorthand for $\mathbf{A} + \mathbf{a} + \mathbf{m}$, where \mathbf{a} is plus and \mathbf{m} is minus) - p stands for "point" - E is a symbol of Σ • Real literals: $sd^+(pd^+ + pd^+Esd^+ + Esd^+)$ • Examples: +6.25, 6.25, -6.25E+2, 6.25E-2, -2E2 Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

ils rineda, ilmas, UNAM & OSU-CIS, 2005

Applications of Regular Expressions

- Provide a "picture" of a pattern that we want to recognize
- They can be "compiled" into determinist automata, which can be modeled to recognize patterns in texts
- Two important applications:
 - Lexical analyzers
 - Texts search in Internet

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

UNIX Notation for RE

- Σ = The set of ASCII characters
- The grep command:
- Global (search for) Regular Expressions and PrintShort hand definitions: Character classes
 - The dot "." stands for any character
 - $[a_1a_2...a_k]$ stands for the *RE*: $a_1 + a_2 + ... + a_k$
 - e.g. The characters used for comparison in C: [<>=!] - [x-y] stand for range definitions:
 - e.g. [A-Za-z0-9] stands for the set of all letters and digits
 A minus sign "---" is placed first or last (to avoid confusion):
 - [-+.0-9] is the set {-, +, . , 0...9} - For reserved characters of UNIX, we use the backslash \
 - [0-9\.] is the set of digits and the dot (not any character) Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 200

UNIX Notation for RE

- Meaning of UNIX operators
 - is used in place of +
 - The operator ? means zero or one of: \mathbf{R} ? = $\Lambda + R$
 - + means one or more: **R**+ = RR^*
 - The operator $\{n\}$ means *n* copies of: **R**{5} = *RRRR*
 - * in UNIX has the usual meaning (not a superscript!)
- Also:
 - [:digit:] stands for [0-9] (not necessarily in ASCII)
 - [:alpha:] stands for [A-Za-z]
 - [:alnum:] stands for [A-Za-z0-9]
- Operators precedence is as usual (with ?, + and $\{n\}$ treated like *)
- UNIX extensions to name and refer to previous strings that have matched a pattern (allowing to the recognition of non-regular languages) are not considered here!

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 200

Lexical Analyzers Lexical analyzer: the part of a compiler that scans the source code and identifies tokens (i.e. basic or atomic.

- source code and identifies *tokens* (i.e. basic or atomic symbols, or entries to the symbol's table)
 - Keywords
 - Identifiers (names, variables, etc.)
- Lexical-analyzer generator: – UNIX's lex (flex in GNU)
 - Accepts a list of regular expressions each followed by a a bracketed piece of code, indicating what to do when an instance of the token described by the *RE* is found
- Advantages:
 - A high level description of a lexical analyzer
 - Automatic generation of complicated code
 - Easy to create and modify
 - Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003

Example

- Partial input to the *lex* command: Else {return
 - Else {return(ELSE)} [A-Za-z][A-Za-z0-9]* {code to Enter identifier in Symbol table;

 - =

{return(GE);} {return(OE):}

return(Id); }

- ...integers, floating-point, character strings, etc.
- Conversion of regular expressions to an automata for
- processing the corresponding strings

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 20

Finding patterns in texts

- *RE* are useful for describing searches for interesting patterns
- Descriptions of vaguely defined class of patterns in texts
- Patters that are hard to define...
- Easy to specify and modify

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003



Example: Detecting addresses in web-pages

- But what about
 - Streets with a different name: "Boulevar", "Place",...
 - Streets with ordinal abbreviations: 42nd St.
 - Post-Office boxes or rural-delivery routes
 - Streets names that do not end with "Street", like *El Camino Real* in Silicon Valley (Spanish name for The Royal Road)
 - El Camino Real Road? • 2000 El Camino Real
- It is really a knowledge engineering task!
- We can appreciate the power of Regular expressions
 - Expressive
 - Economical

Dr. Luis Pineda, IIMAS, UNAM & OSU-CIS, 2003