

RegularExpression

- Java supports a Perl-like regular expression syntax
 - Regular expressions are compiled expressions (patterns) that allow to express syntax rules
 - They enable you to find and replace text an based on an exact syntax (language) (e.g, an email address, creditcard number, zip code etc)

| Expression | example matches | | | | |
|--------------------------------|-----------------|--------|---------|------------|-----------------------|
| <code>mon wed fri</code> | mon | tue | wed | thu | fri sat sun abc 123.2 |
| <code>colo.?r</code> | colour | color | | | |
| <code>[0-9]*Kg</code> | 100gr | 200Kg | 300 | 100MB 12KM | 150Kg |
| <code>[0-9]{2}-[0-9]{3}</code> | 11-891 | 111-90 | 992-172 | a1-222 | |

- Main classes
 - Main class is `Pattern` - compiled regular expression

| Method | | Description |
|---------|--------------------------------------|---|
| Factory | <code>compile (expr)</code> | Creates the <code>Pattern</code> for the specific expressions |
| | <code>compile (expr, flags)</code> | Creates the <code>Pattern</code> using different flags |
| use | <code>split (CharSequence)</code> | Return an array of string particles, split on matches |
| | <code>matcher(CharSequence)</code> | Returns a <code>Matches</code> (used to find or replace text) |

- The `Pattern` is instantiated passing the expression in one of the `compile` methods

```
Pattern p = Pattern.compile("3[47][0-9]{13}");  
... = p.matcher(paymentStr);
```

- An unchecked `PatternSyntaxException` is thrown when pattern is invalid.

```
try{  
    // Invalid expression  
    Pattern pattern =Pattern.compile("10{x}");  
}catch(PatternSyntaxException pse){  
    System.out.println("Invalid regular expression");  
}
```

- `Matcher` class used for matching (find based on different scenarios) and replacing text

- Matching

| Method | | Description |
|------------|---------------------------|--|
| Matching | matches | Attempts to match the entire string |
| | lookingAt | Attempts to match the start of the string |
| | find | Find the next subsequence that matches the pattern |
| | find(int startAt) | Find starting at a specified position |
| match info | start | The start index of the previous match |
| | end | The end index of the previous match |
| | group | Returns the subsequence of previous match |

- Replacing

| Method | Description |
|------------------------------|--|
| replaceAll (String) | Replaces every match with the supplied string |
| replaceFirst (String) | Replace the first match with the supplied string |
| appendReplacement | Append/replace the previous match (append position) |
| appendTail | Append the last part of the input (from append position) |

- Pattern flags

| Flag | Description |
|-------------------------|---|
| CASE_INSENSITIVE | For case-insensitive matching |
| MULTILINE | Position markers are aware of the newline character |
| DOTALL | The "any character" includes the line terminator |
| UNICODE_CASE | Uses unicode instead of ASCII |
| CANON_EQ | Enables canonical equivalence |
| UNIX_LINES | Enable only Unix lines mode |
| LITERAL | The pattern is literal (special characters are literal) |
| COMMENTS | Permits whitespace and comments (are both ignored) |

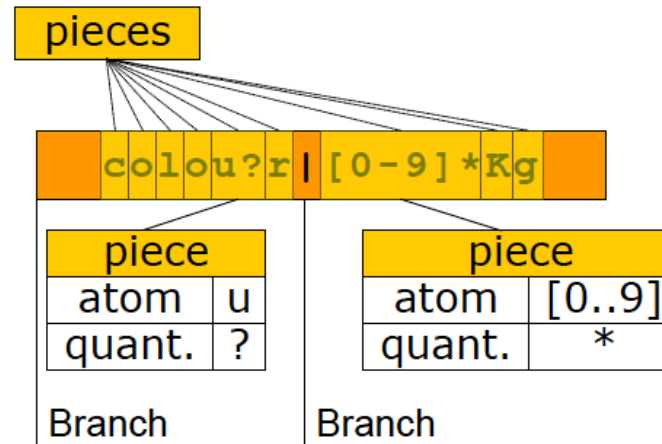
- The `String` class contains methods which take regular expressions

- Without the direct use of `Pattern` and `Matcher` object

```
public boolean matches(String regex)
public String replaceFirst(String regex,String replacement)
public String replaceAll(String regex,String replacement)
public String[] split(String regex,int limit)
public String[] split(String regex)
```

- These methods might throw a `PatternSyntaxException`

- each branch is made up of multiple *pieces* which is an *atom* and an optional *quantifier*



- An atom is a regular expression particle (which has an occurrence)

| Atom type | | | Examples | | | |
|-----------------|------------------|---|--|----------|--|-----------|
| character | | | any character except metacharacters | | | |
| character class | escape class | Predefined characters sequences | <code>\n</code> | newline | <code>.</code> | wildcard |
| | | | <code>\d</code> | digits | <code>\p{M}</code> | marks |
| | expression class | Character groups (negative, positive or intersection) | <code>[abc]</code> | a,b or c | <code>[^ab]</code> | no a or b |
| | | | <code>[a-z]</code> | a to z | <code>\p{Sc}</code> | currency |
| <code>(</code> | <i>regex</i> | <code>)</code> | A group. Contains another regular expression | | <code>([0..9] na)*</code> sequence with digits and 'na' strings | |

- The *metacharacters* are reserved characters that need to be escaped if used literally.

| Metacharacter | |
|---------------|----------------------------|
| . | wild card escape character |
| \ | escape character |
| [] | character class expression |
| () | grouping character |
| ? * + { } | quantifier characters |

| Escaping metacharacter examples | | | |
|---------------------------------|-----|---------|-------|
| 1\.5 | 1.5 | 1\+0\.5 | 1+0.5 |
| 3\\4 | 3\4 | 3\\4 | 3\4 |

- A (greedy) quantifier specifies the occurrence of the atom (quantity is 1 when omitted)

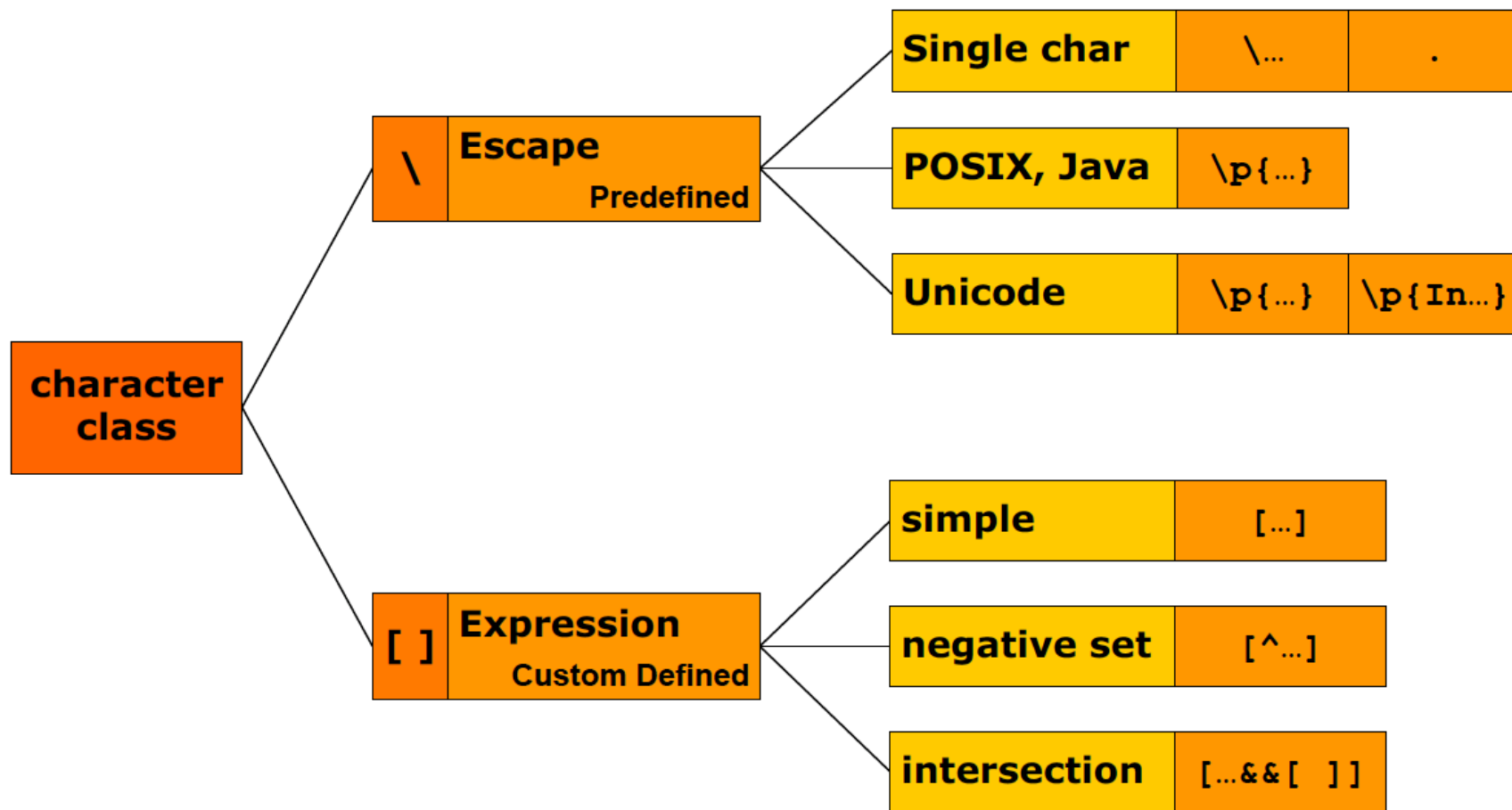
| Quantifier | | # | Example | matches |
|--------------|-------|--------|-------------------|--|
| multiplicity | * | 0..n | How.*\? | "How?", "How are you?", " How much?! " |
| | + | 1..n | http://.+/ | "http://www.w3c.org/", " http://abc " |
| | ? | 0 or 1 | [+-]?[1-9] | "+1", "-9", "4", "-4", "+4" |
| range | {x,y} | range | [0-9]{1,3} | "1", "28", "133", " 9810 ", " 1 ", |
| | {x,} | min | [0-9]{2,} | "10", "299", "392039029", "8" |
| | {x} | exact | [0-9]{2}-[0-9]{2} | "12-32", "00-01", " 1-12 ", " 100-10 " |

- Boundary

| anchor | | example matches |
|-----------------------|------------|--|
| <code>\bex</code> | word | The ex travagant boy ordered T ex -M ex |
| <code>ing\z</code> | end | ing rid, who was listen ing , ordered another thi ing |
| <code>\Aing</code> | start | ing rid, who was listen ing , ordered another thi ing |
| <code>^ing</code> | linestart* | ing rid, who was listen- ing , ordered another thi ing |
| <code>ing&</code> | lineend* | ing rid, who was listen ing , ordered another thi ing |

*Work only in MULTILINE mode (otherwise same as \z and \A)

- *Character classes* define a class of characters either predefined (through escaping) or custom defined by an expression



- Escape characters classes reference predefined character classes

- *Single-character escape* used for character classes with a single character (Metacharacters and other additional characters)

| | | | |
|-----------|------------------------|------------|---------------------|
| \n | Linefeed (\u0A) | \f | form feed (\u0C) |
| \r | carriage return (\u0D) | \a | alert / bell (\u09) |
| \t | tab (\u09) | \e | Escape (\u1B) |
| \- | hyphen | \cX | Control char X |
| \^ | circumflex | \ | pipe |

- *Multi-character escape* used for character classes with multiple characters based on Perl groups, which identify commonly used sets of characters

| Escape sequence | | Explanation |
|-----------------|-----------------|--|
| \s | spaces | All Whitespace characters [#x20\t\n\r] |
| \d | decimal digits | 9, 1, ٢, ٩ (all digits in Unicode) |
| \w | word characters | All except Unicode punctuation, separator and 'other' characters |
| \S \D \W | Complement | Negative set (e.g., \s is all but spaces) |

| \d*Kg (?#Digit followed by Kg)

| \d*\sKg (?#Digit followed a space, followed by Kg)

- The *wildcard* is a special Multi-character escape matching any character (except \n and \r)

| (?#Sentence ending with a question mark)

| .* \?

| How.+ \?

- POSIX character classes (US-ASCII)

| POSIX (\p{...}) | | | |
|--------------------------|--------------|---------------|-------------------|
| Lower | Lower case | Punct | Punctuation |
| Upper | Upper case | Graph | Visible character |
| ASCII | All ASCII | Print | Printable |
| Alpha | Alphabetic | Blank | space or tab |
| Digit | Decimal | Cntrl | Control |
| Alnum | Alphanumeric | XDigit | Hexidecimal |
| | | Space | whitespace |

- *Category escape* used for character classes with multiple characters based on Unicode *general categories* and *blocks*

| Character class | Syntax | Complement | example | |
|------------------|-----------|------------|-------------|-------------------|
| Unicode Category | \p{...} | \P{...} | \p{Lu} | Uppercase letters |
| Unicode Block | \p{In...} | \P{In...} | \p{InGreek} | Greek characters |

- Unicode *categories* and *blocks*

- Unicode categories

| Letters (L) | | | Numbers (N) | | | Symbols (S) | | | Punctuation (P) | | |
|-------------|------------|----|----------------|----------|-------|-------------|----------|---|-----------------|---------------|---|
| Lu | Uppercase | Δ | Nd | Decimal | 1 | Sm | Math | + | Pc | Connector | — |
| Ll | Lowercase | a | NI | Number | xi | Sc | Currency | € | Pd | Dash | - |
| Lt | Title case | Lj | No | Other | 2 | Sk | Modifier | ¨ | Ps | Open | (|
| Lm | Modifiers | ^ | Separators (Z) | | | So | Other | © | Pe | Close |) |
| Lo | Other | 2 | | | | | | | Pi | Initial quote | " |
| Not shown | | | Zs | Space | x20 | | | | Pf | Final quote | " |
| | | | Zl | Line | x2028 | | | | Po | Other | ! |
| | | | Zp | Pargrph. | x2029 | | | | | | |

Marks (M): Mn (nonspacing), Mc (spacing combining), Me (enclosing)
 Other (C) :Cc (control), Cf (format), Co (private), Cn (not assigned)

(?# Question starting with a uppercase letter)
\p{Lu}.*\?

- Unicode blocks

| Example Unicode Blocks | | | |
|------------------------|------------|--------------------|-------------------|
| BasicLatin | Syriac | Tibetan | Dingbats |
| Latin-1Supplement | Thaana | Runic | BraillePatterns |
| LatinExtended-A | Devanagari | Khmer | YiSyllables |
| LatinExtended-B | Bengali | Mongolian | YiRadicals |
| Greek | Gurmukhi | LatinExtendedAd... | HangulSyllables |
| Cyrillic | Gujarati | GreekExtended | MathematicalOp... |
| Armenian | Tamil | GeneralPunctuation | Arrows |
| Hebrew | Telugu | CurrencySymbols | NumberForms |
| Arabic | Thai | GeometricShapes | ... |

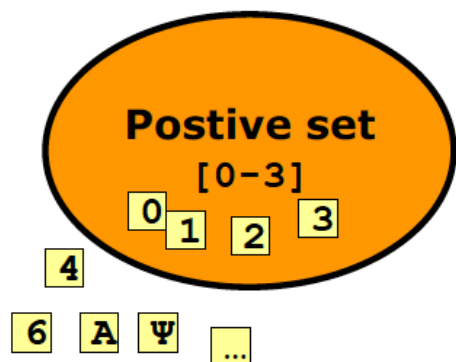
(?# Sequence of basic latin characters)
\p{inBasicLatin}*

- Character class expressions* define custom character classes

positive character groups

Union of atoms

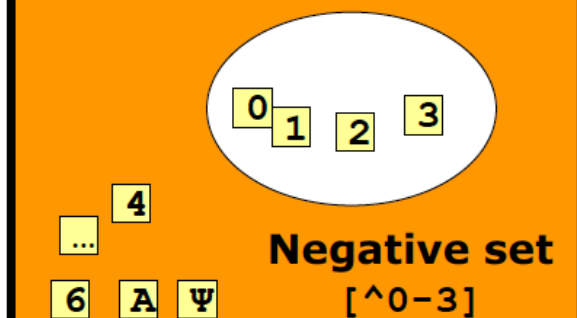
Unicode



negative character group

Absolute complement

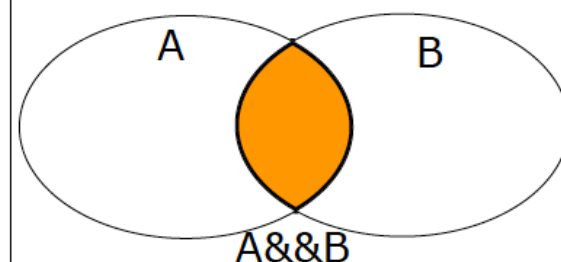
Unicode



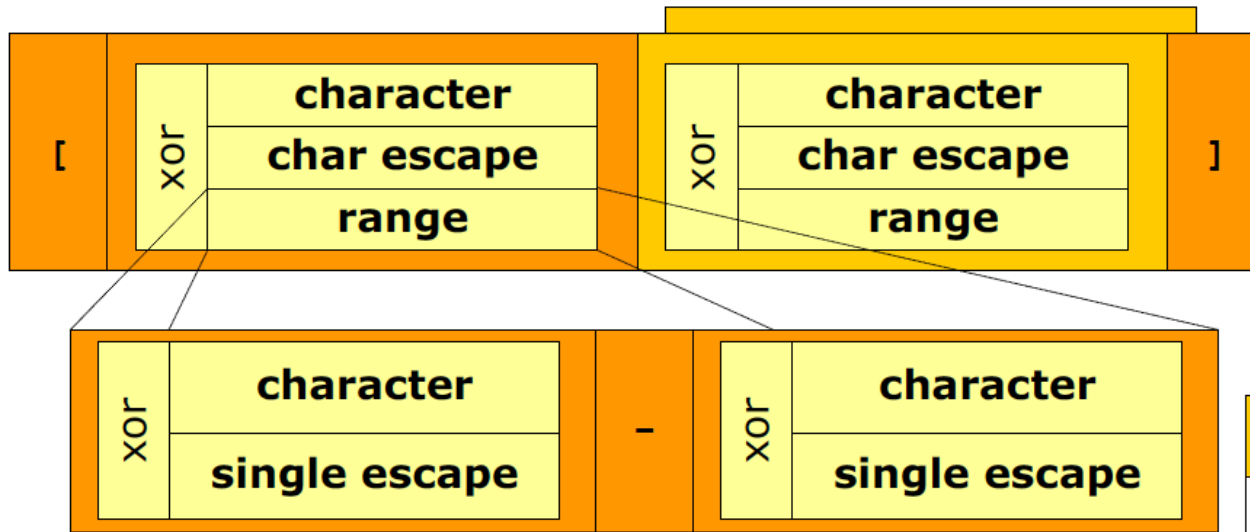
character class intersection

Intersection

Unicode



- The *simple character group* is the basic expression that defines the set of characters as a union of atoms.



is a range of Unicode code points

| Examples | |
|------------|-------------|
| [a] | [0-46-9] |
| [! ?] | [0-4\p{L1}] |
| [a\p{N}\s] | [!-\.] |
| [0-9] | [a-z] |

- The *negative character group* defines an **absolute complement**

| | | | |
|---|---|------------------------|---|
| [| ^ | simple character group |] |
|---|---|------------------------|---|

| Examples |
|-------------|
| [^a] |
| [^! ?] |
| [^a\p{N}\s] |
| [^0-9] |

- The *intersection character group* defines an **intersection**

| Examples | Results in Set |
|--|-----------------------------|
| <code>[\p{inBasicLatin}&[\p{Lu}]]</code> | All Latin uppercase letters |
| <code>[\p{InArabic}&[\p{N}]]</code> | All Arabic numbers |

- Last two can be used to make a subtraction of two sets (**relative complement**)

| | | | | | | | |
|---|--------------------------------------|----|---|---|-----------------|---|---|
| [| positive/negative character group | && | [| ^ | character group |] |] |
|---|--------------------------------------|----|---|---|-----------------|---|---|

| Examples |
|------------------|
| [0-9&&[^3]] |
| [\p{Lu}&&[^M-Z]] |

Questions

