



Regular Expressions Κανονικές εκφράσεις

Αθ. Ανδρούτσος



Κανονικές Εκφράσεις (1)

- Με τις κανονικές εκφράσεις (regular expressions) μπορούμε να ελέγχουμε **εάν μια συμβολοσειρά είναι ίδια με κάποια άλλη ή περιέχει ένα συγκεκριμένο μοτίβο (pattern)**
- Τα regular expressions συχνά χρησιμοποιούνται για αναζήτηση σε log files, για validation των δεδομένων εισόδου ή για την υλοποίηση scanners σε μεταγλωττιστές
- Οι κλάσεις στην Java που υλοποιούν κανονικές εκφράσεις βρίσκονται στο package `java.util.regex`



String.matches

- Μία απλή εκδοχή για τη χρήση Regular Expressions είναι η μέθοδος ***String.matches(String regex)***
- Είναι παρόμοια με την ***equals*** μόνο που εδώ συγκρίνουμε με regular expression και επομένως συγκρίνουμε με πολλαπλές τιμές
- Έστω για παράδειγμα ότι θέλουμε να ελέγξουμε αν ένα String περιέχει τη λέξη Red (βλ. επόμενη διαφάνεια)



Απλή εφαρμογή της .matches

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * Περιέχει μία μέθοδο που ελέγχει
5  * αν ένα String είναι το "Red".
6  *
7  * @author A. Androutsos
8  * @see #isRed(String)
9  */
10 public class RegExEquals {
11
12     @ public static boolean isRed(String s) {
13         return (s != null) && (s.matches("Red"));
14     }
15 }
```

- Εδώ η .matches λειτουργεί όπως η equals
- Ελέγχει απλώς αν το String s είναι "Red"

- Αν θα θέλαμε να ελέγξουμε αν το String s είναι "Red" ή "Green" ;
- Με την .equals θα είχαμε: s.equals("Red") || s.equals("Green")
- Με Regular Expression δείτε την επόμενη διαφάνεια



Alternation -- Red or Green

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * Περιέχει μία μέθοδο που ελέγχει αν
5  * ένα String είναι το "Red" ή το "Green".
6  *
7  * @author A. Androutsos
8  * @see #isRedOrGreen(String)
9  */
10 public class RegExRedOrGreen {
11
12     public static boolean isRedOrGreen(String s) {
13         return (s != null) && (s.matches("Red|Green"));
14     }
15 }
```

- Παρατηρήστε ότι **δεν αφήνουμε κενά μεταξύ των λέξεων και του /**

- Αν θέλουμε να μπορεί να ταιριάζει το String s με τα Red ή red ή Green ή green ;
- Δείτε την επόμενη διαφάνεια.



Character Sets

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * Περιέχει μία μέθοδο που ελέγχει αν
5  * ένα String είναι το "Red" ή "red" ή
6  * "Green" ή "green".
7  *
8  * @author A. Androutsos
9  * @see #isRredOrGgreen(String)
10 */
11 public class RegExRedOrGreenFirstLetter {
12
13     public static boolean isRredOrGgreen(String s) {
14         return (s != null) && (s.matches("[Rr]ed|[Gg]reen"));
15     }
16 }
```

- Έλεγχος για Red/red/Green/green
- Μέσα σε [] δίνουμε τους χαρακτήρες που μπορεί να βρίσκονται στην 1^η θέση του Red ή του Green

- Αν θα θέλαμε μία λέξη να ξεκινάει με οποιοδήποτε κεφαλαίο γράμμα ;
- Δείτε την επόμενη διαφάνεια



Έλεγχος περιοχής χαρακτήρων (1)

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * Περιέχει μία μέθοδο που ελέγχει αν
5  * ένα String ξεκινάει με ένα γράμμα κεφαλαίο
6  * και τελειώνει με -ing.
7  *
8  * @author A. Androutsos
9  * @see #hasOneCapitalLetterFirstAndEndsWith_ing(String)
10 */
11 public class RegExRange {
12
13     public static boolean hasOneCapitalLetterFirstAndEndsWith_ing(String s) {
14         return (s != null) && (s.matches("[A-Z]ing"));
15     }
16 }
```

Με παύλα
ορίζουμε
την περιοχή

Για
παράδειγμα
με [A-Z]
ορίζουμε
ένα
κεφαλαίο
χαρακτήρα

- Αν θα θέλαμε να ξεκινάει είτε με κεφαλαίο ή με μικρό γράμμα ;
- Δείτε την επόμενη διαφάνεια



Έλεγχος περιοχής χαρακτήρων (1)

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * Περιέχει μία μέθοδο που ελέγχει αν
5  * ένα String ξεκινάει με ένα γράμμα πεζό
6  * ή κεφαλαίο και τελειώνει με -ing.
7  *
8  * @author A. Androutsos
9  * @see #hasOneCapitalLetterFirstAndEndsWith_ing(String)
10 */
11 public class RegExRange2 {
12
13     public static boolean hasOneCapitalLetterFirstAndEndsWith_ing(String s) {
14         return (s != null) && (s.matches("[a-zA-Z]ing"));
15     }
16 }
```

- Με παύλα ορίζουμε την περιοχή, π.χ. με **[a-zA-Z]** ορίζουμε ένα είτε πεζό ή κεφαλαίο χαρακτήρα



Μεταχαρακτήρες (1)

Κανονικές Εκφράσεις

- . Οποιοσδήποτε χαρακτήρας
- ^ Αρχή μιας γραμμής
- \$ Τέλος μιας γραμμής
- \w Χαρακτήρας λέξης από a-z, A-Z, 0-9, καθώς και _ (underscore)
- \W Οποιοσδήποτε χαρακτήρας εκτός από χαρακτήρα λέξης



Μεταχαρακτήρες (2)

Κανονικές Εκφράσεις

- **\s** Χαρακτήρας διαστήματος
- **\S** Οποιοσδήποτε χαρακτήρας εκτός από χαρακτήρα διαστήματος
- **\d** Οποιοδήποτε ψηφίο
- **\D** Οποιοσδήποτε χαρακτήρας εκτός από ψηφίο



Οποιοσδήποτε χαρακτήρας - Κενό - Ψηφίο

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * Περιέχει μία μέθοδο που ελέγχει αν
5  * ένα String έχει τη μορφή ένας οποιοσδήποτε
6  * χαρακτήρας ακολουθούμενος από ένα κενό διάστημα
7  * ακολουθούμενος από ένα ψηφίο.
8  *
9  * @author A. Androutsos
10  * @see #isAnySpaceDigit(String)
11  */
12 public class RegExAnyCharSpaceDigit {
13
14     @
15     public static boolean isAnySpaceDigit(String s) {
16         return (s != null) && (s.matches("\\s\\d"));
17     }
18 }
```

- Θέλουμε ως 1ο χαρακτήρα οποιοδήποτε χαρακτήρα. Για αυτό χρησιμοποιούμε τον μεταχαρακτήρα τελεία.
- Στη συνέχεια θέλουμε ένα κενό (space \s) και μετά ένα ψηφίο (digit \d)
- Ως ψηφίο θα μπορούσαμε να χρησιμοποιήσουμε και το [0-9]

- Παρατηρήστε πως επειδή το \ είναι ειδικός χαρακτήρας τον κάνουμε escape με \ μπροστά, οπότε έχουμε \\s για το κενό διάστημα που είναι \s και \\d για το ψηφίο που είναι \d



Γράμμα - Ψηφίο

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * Περιέχει μία μέθοδο που ελέγχει αν
5  * ένα String έχει τη μορφή ένα γράμμα
6  * ή ψηφίο ή κάτω παύλα ακολουθούμενο
7  * από ένα ψηφίο.
8  *
9  * @author A. Androutsos
10 * @see #isLetterDigit(String)
11 */
12 public class RegExLetterDigit {
13
14 @ public static boolean isLetterDigit(String s) {
15     return (s != null) && (s.matches("\\w\\d"));
16 }
17 }
```

- Ένα οποιοδήποτε γράμμα δηλαδή [a-zA-Z0-9_] αναπαρίσταται ως \w και ένα ψηφίο [0-9] ως \d
- Εναλλακτικά θα μπορούσαμε να το γράψουμε στην πλήρη του μορφή



Γράμμα και ψηφίο οπουδήποτε μέσα στο String

Κανονικές Εκφράσεις

- Η διαφορά με την προηγούμενη διαφάνεια είναι ότι έχουμε προσθέσει στην αρχή και το τέλος το `(.*)` που σημαίνει ένας οποιοσδήποτε χαρακτήρας μηδέν ή περισσότερες φορές

Το `*` είναι ποσοδείκτης που σημαίνει μηδέν ή περισσότερες φορές αυτό που προηγείται, εν προκειμένω η τελεία που σημαίνει οποιοσδήποτε χαρακτήρας
Strings που κάνουν match είναι για παράδειγμα τα: ***F4*** , ***ABCD889!*** , ***55AA***

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * @author A. Androutsos
5  * @see #isLetterDigitAnywhere(String)
6  */
7 public class RegExLetterDigitAnywhere {
8
9     /**
10      * Ελέγχει αν ένα String ξεκινάει με ένα οποιοδήποτε
11      * χαρακτήρα θ ή περισσότερες φορές βρίσκει ένα χαρακτήρα
12      * λέξης και στη συνέχεια ένα ψηφίο και τελειώνει με ένα
13      * οποιοδήποτε χαρακτήρα θ ή περισσότερες φορές.
14      *
15      * @param s το προς έλεγχο String
16      * @return true, αν γίνει match το regex, αλλιώς false.
17      */
18     @ public static boolean isLetterDigitAnywhere(String s) {
19         return (s != null) && (s.matches(".*\\w\\d.*"));
20     }
21
22     /**
23      * Είναι παρόμοιο με:
24      * return (s != null) && (s.matches(".*[a-zA-Z0-9_][0-9].*"));
25      */
26 }
```



Γράμμα και ψηφίο ως λέξη

Κανονικές Εκφράσεις

- Η διαφορά με την προηγούμενη διαφάνεια είναι ότι έχουμε προσθέσει το `\b` στην αρχή και το τέλος που σημαίνει word boundary άρα **αναζητούμε λέξεις** που αποτελούνται από ένα γράμμα και ένα ψηφίο

Strings που κάνουν match είναι για παράδειγμα τα:
P **F4** , ABC **D8** 89! ,
55 **A0**

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * @author A. Androutsos
5  * @see #isLetterDigitAsWholeWord(String)
6  */
7 public class RegExLetterDigitAsWord {
8
9     /**
10      * Ελέγχει αν ένα String ξεκινάει με ένα οποιοδήποτε
11      * χαρακτήρα θ ή περισσότερες φορές βρίσκει ένα χαρακτήρα
12      * word boundary, μετά ένα χαρακτήρα λέξης, μετά ένα ψηφίο,
13      * μετά ένα χαρακτήρα word boundary (άρα βρίσκει λέξεις)
14      * και τελειώνει με ένα οποιοδήποτε χαρακτήρα θ ή περισσότερες
15      * φορές.
16      *
17      * @param s το προς έλεγχο String
18      * @return true, αν γίνει match το regex, αλλιώς false.
19      */
20     @ public static boolean isLetterDigitAsWholeWord(String s) {
21         return (s != null) && (s.matches(".*\\b\\w\\d\\b.*"));
22     }
```



Ποσοδείκτες

- Όπως το * που είδαμε στο προηγούμενο παράδειγμα, οι ποσοδείκτες εφαρμόζονται αμέσως μετά από τον χαρακτήρα, ή το regex (ή το group)
 - ? Μηδέν ή μία εμφάνιση
 - + Μία ή περισσότερες εμφανίσεις (occurrence indicator)
 - * 0 ή περισσότερες εμφανίσεις
 - {n} n εμφανίσεις
 - {n, m} Μεταξύ n και m εμφανίσεων
 - {n, } Τουλάχιστον n εμφανίσεις
 - {,m} Όχι περισσότερες από m εμφανίσεις



Έγκυρη διεύθυνση e-mail

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * Περιέχει μία μέθοδο που ελέγχει αν
5  * ένα String αναπαριστά μία έγκυρη email
6  * διεύθυνση που τελειώνει σε com ή gr.
7  *
8  * @author A. Androutsos
9  * @see #isValidEmail(String)
10 */
11 public class RegExEmail {
12
13     /**
14      * Ελέγχει αν ένα String αναπαριστά μία έγκυρη email
15      * διεύθυνση.
16      *
17      * @param s το προς έλεγχο String
18      * @return true, αν γίνει match το regex, αλλιώς false.
19      */
20     @ public static boolean isValidEmail(String s) {
21         return (s != null) && (s.matches("\\w*\\.?\\w+@\\w+\\.?(com|gr)"));
22     }
23 }
```

- Ελέγχει αν το s περιέχει έγκυρη διεύθυνση e-mail που τελειώνει σε com ή gr και έχει τη μορφή 0 ή περισσότεροι χαρακτήρες, 0 ή μία τελεία –επειδή η τελεία είναι ειδικός χαρακτήρας για να πάρουμε την τελεία την κάνουμε escape \\.
- Στη συνέχεια πάλι χαρακτήρες ένα ή περισσότερους, μετά το @, μετά πάλι ένα ή περισσότερους χαρακτήρες , μετά τελεία και μετά com ή gr

- Το com|gr το βάζουμε σε παρενθέσεις επειδή πρόκειται για ξεχωριστό regex group



Επαλήθευση format

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * @author A. Androutsos
5  * @see #isTwoLettersDashSevenDigits(String)
6  */
7 public class RegExTwoLettersDashSevenDigits {
8
9     /**
10      * Ελέγχει αν ένα String ξεκινάει με δύο γράμματα
11      * κεφαλαία στη συνέχεια έχει παύλα και μετά 7 ψηφία.
12      *
13      * @param s το προς έλεγχο String
14      * @return true, αν γίνει match το regex, αλλιώς false.
15      */
16     @ public static boolean isTwoLettersDashSevenDigits(String s) {
17         return (s != null) && (s.matches("[A-Z]{2}-\\d{7}"));
18     }
19 }
```

- Κάνουμε validate ένα string ότι περιέχει τη μορφή δύο γράμματα κεφαλαία ακολουθούμενα από παύλα – ακολουθούμενα από 7 αριθμούς



Split με RegEx

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * Επιδεικνύει την String.split με
5  * Regular Expression.
6  *
7  * @author A. Androutsos
8  */
9 public class RegExSplit {
10
11     /**
12      * Κάνει parse το String και το επιστρέφει
13      * τα tokens που διαχωρίζονται με ένα ή
14      * περισσότερα κενά διαστήματα.
15      */
16     public static void main(String[] args) {
17         String s = "Java    Advanced Level";
18
19         String[] tokens = s.split("\\s+");
20
21         for (String token : tokens) {
22             System.out.println(token);
23         }
24     }
25 }
```

- Στην `split` δίνουμε το regex `\\s+` που σημαίνει το να διαχωρίσει το string σε tokens που με βάση ένα ή περισσότερα κενά διαστήματα



replaceAll με RegEx

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 /**
4  * Επιδεικνύει την String.split με
5  * Regular Expression.
6  *
7  * @author A. Androutsos
8  */
9 public class RegExReplaceAll {
10
11     /**
12      * Κάνει parse το String και το επιστρέφει
13      * τα tokens που διαχωρίζονται με ένα ή
14      * περισσότερα κενά διαστήματα.
15      */
16     public static void main(String[] args) {
17         String s = "Java    Advanced Level";
18
19         s = s.replaceAll("\\s+", " ");
20         System.out.println(s);
21     }
22 }
```

"C:\Program Files\Ja
Java Advanced Level
Process finished wit

- Αντικαθιστούμε τα κενά με ένα κενό



Αμοιβαία ανταλλαγή με backreference

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 public class RegExBackReference {
4
5     public static void main(String[] args) {
6
7         String s = "Thanassis Androutsos";
8
9         s = s.replaceAll("(.*)(\\s)(.*)", "$2 $1");
10
11         System.out.println(s);
12     }
13 }
```

Execution output: "C:\Program Files\Java\jdk-11.0.2\bin\java.exe" -cp "C:\Program Files\Java\jdk-11.0.2\bin\java.exe" gr.aueb.elearn.ch3.RegExBackReference
Androutsos Thanassis
Process finished with exit code 0

- Groups ορίζουμε μέσα σε παρενθέσεις, ώστε να μπορούμε να αναφερθούμε σε αυτά αργότερα (backreference)
- Τα \$1 , \$2 αναφέρονται στα δύο groups



Pattern - Matches

- Με την `matches` δεν μπορούμε να χειριστούμε πολλαπλές εμφανίσεις ενός pattern μέσα σε ένα string
- Χρησιμοποιούμε τις κλάσεις `Pattern` / `Matcher`
- Η κλάση `Pattern` μας δίνει τη δυνατότητα να κάνουμε `compile` ένα regex μία φορά και να το χρησιμοποιήσουμε όσες φορές θέλουμε και επομένως είναι πολύ αποδοτικότερη



Μετατροπή Ημερομηνίας

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 /**
7  * Μετατρέπει ημερομηνίες από mm/dd/yyyy
8  * σε dd/mm/yyyy.
9  *
10 * @author A. Androutsos
11 */
12 public class PatternDate {
13
14     public static void main(String[] args) {
15         String s = "12/05/1996";
16
17         Pattern pattern = Pattern.compile("(\\d{2})/(\\d{2})/(\\d{4})");
18         Matcher matcher = pattern.matcher(s);
19
20         //Attempts to match the entire region against the pattern.
21         if (matcher.matches()) {
22             String month = matcher.group(1);
23             String day = matcher.group(2);
24             String year = matcher.group(3);
25             String out = day + "/" + month + "/" + year;
26             System.out.println(out);
27         } else {
28             System.out.println("Pattern not matches");
29         }
30     }
31 }
```

- Μέσα σε παρενθέσεις είναι τα groups ώστε στη συνέχεια να τα επεξεργαστούμε



Greedy Ποσοδείκτης

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 public class PatternGreedy {
7
8     public static void main(String[] args) {
9         String s = "login:thanos;pass:123456;";
10
11         Pattern pattern = Pattern.compile(".*;");
12
13         Matcher matcher = pattern.matcher(s);
14
15         while (matcher.find()) {
16             System.out.println(matcher.group(0));
17         }
18     }
19 }
```

"C:\Program Files\Java\jdk-
login:thanos;pass:123456;

Process finished with exit

- Το match γίνεται μέχρι το τελευταίο ;
Δηλ. As much as possible (greedy – άπληστο match)
- Η find() δίνει το next match

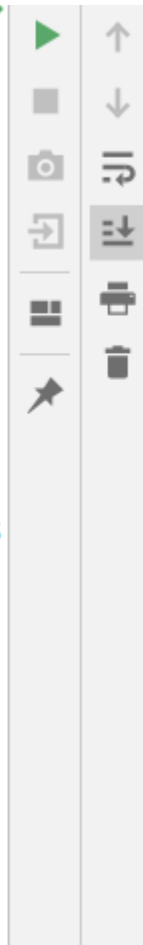
- Στο pattern δεν έχουμε ορίσει groups, οπότε όλο το match θεωρείται group(0) ή σκέτο group()



Reluctant (non-greedy) ποσοδείκτης

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 public class PatternReluctant {
7
8     public static void main(String[] args) {
9         String s = "login:thanos;pass:123456;";
10
11         Pattern pattern = Pattern.compile(".*?;");
12
13         Matcher matcher = pattern.matcher(s);
14
15         while (matcher.find()) {
16             System.out.println(matcher.group(0));
17         }
18     }
19 }
```



"C:\Program Fi

login:thanos;

pass:123456;

Process finish

- Το match γίνεται μέχρι το πρώτο ;
Δηλαδή
As little as possible (non-greedy – όχι άπληστο match)



Start of String

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 public class PatternStartOfString {
7
8     public static void main(String[] args) {
9         String s = "B31:F456:D55";
10
11         Pattern pattern = Pattern.compile("^([A-Z].*?):");
12         Matcher matcher = pattern.matcher(s);
13
14         while (matcher.find()) {
15             System.out.println(matcher.group());
16         }
17     }
18 }
```

- Με το `^` έξω από τα `[]` εννοούμε την αρχή του string
- Επομένως ψάχνουμε για pattern που ξεκινάει από την αρχή με A-Z και φτάνει μέχρι τα πρώτα :



Αφαίρεση συνόλου

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 public class PatternNegation {
7
8     public static void main(String[] args) {
9         String s = "A31:D456:X55:H78:";
10
11         Pattern pattern = Pattern.compile("[A-Z&&[^ADX]].*?:");
12         Matcher matcher = pattern.matcher(s);
13
14         while (matcher.find()) {
15             System.out.println(matcher.group());
16         }
17     }
18 }
```

"C:\Pr
H78:

Δίνει H78: μιας
και τα
υπόλοιπα
substrings
εμπίπτουν στο
^ADX

Με το ^ μέσα
στα []
εννοούμε
άρνηση
(negation)

- Το `[A-Z&&[^ADF]]` δίνει την αφαίρεση από το A-Z του `^ADX`



Τομή συνόλων

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 public class RegExIntersection {
7
8     public static void main(String[] args) {
9         String s = "A31:D456:X55:H78:";
10
11         Pattern pattern = Pattern.compile("[A-Z&&[ADX]].*?:");
12         Matcher matcher = pattern.matcher(s);
13
14         while (matcher.find()) {
15             System.out.println(matcher.group());
16         }
17     }
18 }
```

"C:\Pro
A31:
D456:
X55:

Proces:

Δίνει τα παραπάνω μιας και το H78: δεν εμπίπτει στο ADX

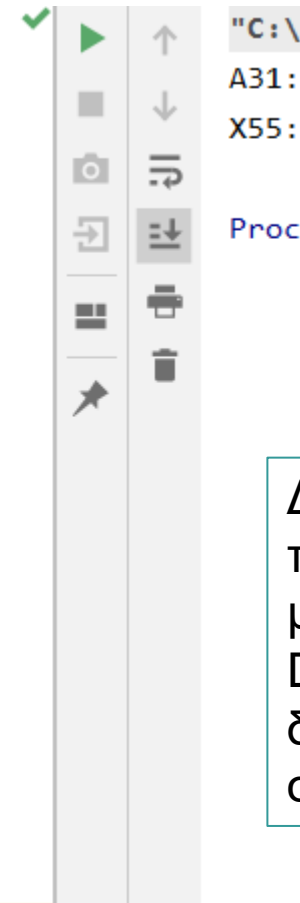
- Το `[A-Z&&[ADF]]` δίνει την τομή του A-Z και του ADX, άρα δίνει τα ADX



Ένωση Συνόλων

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 public class RegExUnion {
7
8     public static void main(String[] args) {
9         String s = "A31:D456:X55:H78:";
10
11         Pattern pattern = Pattern.compile("[A-C[M-Z]].*?:");
12         Matcher matcher = pattern.matcher(s);
13
14         while (matcher.find()) {
15             System.out.println(matcher.group());
16         }
17     }
18 }
```



Δίνει τα
παραπάνω
μιας και τα
D465, H78:
δεν εμπίπτουν
στο [A-C[M-Z]]

- Το [A-Z[M-Z]] δίνει την ένωση του A-C και του M-Z



Groups & Capturing

Κανονικές Εκφράσεις

- Όπως είδαμε στα προηγούμενα παραδείγματα, βάζοντας μία έκφραση μέσα σε παρενθέσεις δημιουργούμε capturing groups δηλ. named groups στα οποία μπορούμε στη συνέχεια να αναφερθούμε (backreference), όπως στο παράδειγμα με τις ημερομηνίες
- Ο τρόπος αναφοράς είναι με αριθμούς ξεκινώντας από το 1 - group(1)



Named Groups

- Σε ένα group μπορούμε να δώσουμε και ένα όνομα. Όταν δώσουμε όνομα σε ένα group μπορούμε στη συνέχεια να αναφερθούμε στο group με το όνομά του αλλά εξακολουθητικά και με την αρίθμηση που αποδίδεται αυτόματα

```
1 package gr.aueb.elearn.ch3;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 /**
7  * Μετατρέπει ημερομηνίες από mm/dd/yyyy
8  * σε dd/mm/yyyy με τη χρήση named groups.
9  *
10 * @author A. Androutsos
11 */
```

Δείτε τη
συνέχεια στην
επόμενη
διαφάνεια



Named Groups

Κανονικές Εκφράσεις

```
12 ▶ public class PatternDateNamed {
13
14 ▶ public static void main(String[] args) {
15     String s = "12/05/1996";
16
17     Pattern pattern = Pattern.compile("(?<month>\\d{2})/(?<day>\\d{2})/(?<year>\\d{4})");
18     Matcher matcher = pattern.matcher(s);
19
20     //Attempts to match the entire region against the pattern.
21     if (matcher.matches()) {
22         String month = matcher.group("month");
23         String day = matcher.group("day");
24         String year = matcher.group("year");
25         String out = day + "/" + month + "/" + year;
26         System.out.println(out);
27     } else {
28         System.out.println("Pattern not matches");
29     }
30 }
31 }
```

Named groups

Αναφορές στα
named groups
Backreferences

- Παρατηρήστε τα named groups και στη συνέχεια τα backreferences με **`matcher.group("groupName")`**



Capturing vs non-capturing groups

Κανονικές Εκφράσεις

- Τα groups που είδαμε με τον τρόπο που τα ορίζουμε μέσα σε παρενθέσεις είναι capturing groups δηλ. καταναλώνουν τους χαρακτήρες και μπορούμε να αναφερθούμε σε αυτά, άρα αποτελούν μέρους του αποτελέσματος του regex
- Υπάρχουν όμως περιπτώσεις που θα θέλαμε μόνο να ελέγχουμε την ύπαρξη ενός regex group χωρίς το group να καταναλώσει χαρακτήρες, αλλά απλώς να μας πει αν υπάρχει ή όχι αυτό το pattern στο string. Αυτά τα groups ονομάζονται **non-capturing** και **δεν μπορούμε να αναφερθούμε σε αυτά με backreference**



Non-capturing groups

Κανονικές Εκφράσεις

- Μπορούμε να δηλώσουμε ένα non-capturing group ως εξής:
- `(Hello)(?:World)+`
- Αν έχουμε τα String HelloWorld, τότε γίνεται match αλλά αποδίδεται στο match μόνο το Hello



Capturing groups - HelloWorld

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 public class RegExNonCapturingHelloWorld {
7
8     public static void main(String[] args) {
9         String s = "HelloWorld HelloWorldWorld";
10
11         Pattern pattern = Pattern.compile("(Hello)(World)+");
12         Matcher matcher = pattern.matcher(s);
13
14         while (matcher.find()) {
15             for (int i=1; i <= matcher.groupCount(); i++) {
16                 System.out.print(matcher.group(i) + " ");
17             }
18             System.out.println();
19         }
20     }
21 }
```

```
"C:\Program I
Hello World
Hello World

Process finish
```

- Σε κάθε match (εδώ έχουμε 2 match) έχουμε groupCount groups (εδώ έχουμε δύο groups)

- Όπως βλέπετε γίνεται capture τόσο το HelloWorld, όσο και το HelloWorldWorld, και αποδίδεται το Hello World και στις δύο περιπτώσεις



Non-capturing groups - HelloWorld

Κανονικές Εκφράσεις

```
1 package gr.aueb.elearn.ch3;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 public class RegExNonCapturingHelloWorld {
7
8     public static void main(String[] args) {
9         String s= "HelloWorld HelloWorldWorld";
10
11         Pattern pattern = Pattern.compile("(Hello)(?:World)+");
12         Matcher matcher = pattern.matcher(s);
13
14         while (matcher.find()) {
15             for (int i=1; i <= matcher.groupCount(); i++)
16                 System.out.println(matcher.group(i));
17         }
18     }
19 }
```

- Το `?:` κάνει match το Hello World αλλά δεν κάνει capture to World (κάνει το group non-capturing)

- Όπως βλέπετε ενώ γίνεται capture τόσο το HelloWorld, όσο και το HelloWorldWorld, **αποδίδεται μόνο το Hello** και στις δύο περιπτώσεις



Non-capturing groups

- Τα non-capturing groups λειτουργούν ως **zero length assertions** δηλαδή μπορεί να κάνουν match αλλά δεν καταναλώνουν του χαρακτήρες και δεν αποδίδουν το match (zero-length match)
- Αν θέλουμε να δηλώσουμε ότι ένα regex **ακολουθείται από ένα non-capturing group το οποίο δεν κάνει match** χρησιμοποιούμε το **?**
- Το **?:** κάνει match αλλά δεν κάνει capture
- Το **?=** δεν κάνει ούτε match ούτε capture
- Για παράδειγμα το **^(?=.*?[a-z]).{6,}\$** σημαίνει ένα string που περιέχει τουλάχιστον 6 οποιουσδήποτε χαρακτήρες. Το group είναι non-capturing και κάνει assert ότι το string περιέχει τουλάχιστον ένα μικρό γράμμα
- Το **?** πριν το **[a-z]** κάνει το regex reluctant και βελτιστοποιεί την αναζήτηση



Παράδειγμα

Κανονικές Εκφράσεις

```
1 package gr.aueb.cf.ch20.regexapp;
2
3 import java.util.regex.Matcher;
4 import java.util.regex.Pattern;
5
6 /**
7  * Gets a string containing 8 or more chars and asserts that the string includes
8  * at least one lowercase letter, one uppercase letter and one digit.
9  */
10 public class LowercaseValidationApp {
11
12     public static void main(String[] args) {
13
14         String s = "0Sgwer6ty";
15
16         Pattern pattern = Pattern.compile("(?=.*?[a-z])(?=.*?[A-Z])(?=.*?[0-9])^.{8,}$");
17         Matcher matcher = pattern.matcher(s);
18
19         while (matcher.find()) {
20             System.out.println(matcher.group());
21         }
22     }
23 }
```

- Ελέγχει αν το string περιέχει ένα τουλάχιστον πεζό γράμμα, ένα τουλάχιστον κεφαλαίο και ένα τουλάχιστον ψηφίο



Μικρή εργασία

Κανονικές Εκφράσεις

- Έστω ότι ο χρήστης δίνει ένα password.
- Επαληθεύστε ότι έχει τη μορφή:
 - Τουλάχιστον 8 χαρακτήρες
 - Τουλάχιστον 1 μικρό γράμμα
 - Τουλάχιστον 1 κεφαλαίο γράμμα
 - Τουλάχιστον 1 αριθμό
 - Τουλάχιστον 1 ειδικό χαρακτήρα `#?!@$%^&*-`