

Linguistica Computazionale

Salvatore Sorce
Dipartimento di Ingegneria
Chimica, Gestionale, Informatica e Meccanica

Ludici Adattati da Alessandro Lenci
Dipartimento di Linguistica "T. Bolelli"

Espressioni Regolari

Espressioni Regolari (RE)

“A language for specifying text search strings

A string is a sequence of symbols

- Broadly speaking, any sequence of alphanumeric characters (letters, numbers, spaces, tabs, and punctuation)*

Formally, a Regular Expression is an algebraic notation for characterizing a set of strings.”

Wildcard

- La RE ***/./*** corrisponde a qualsiasi carattere (eccetto il ritorno-a-capo)
- ***/./ /.* /***

RE	Definizione	Esempi di "matching"
<i>/b.s/</i>	qualsiasi stringa di tre caratteri che inizia con 'b' e termina con 's'	<i>"<u>bas</u>" "<u>bbs</u>" "<u>b3s</u>" "<u>b!s</u>" "<u>b s</u>" "<u>b,s</u>"</i> ... <i>"baas"</i>
<i>/b.*s/</i>	qualsiasi stringa che inizia per b e termina per s	<i>"<u>bs</u>" "<u>bas</u>" "<u>bbs</u>" "<u>b3s</u>" "<u>b!s</u>" "<u>b s</u>"</i> <i>"<u>b,s</u>" "<u>baas</u>"</i> <i>"<u>bisogna prendere l'autobus</u>"</i>
<i>./.* /</i>	qualsiasi stringa (compresa quella vuota)	<i>l'intero corpus, anche se vuoto</i>

Raggruppamento e memoria

- Le parentesi tonde servono per **raggruppare** stringhe di caratteri da moltiplicare:

RE	Definizione	Esempi di "matching"
<code>/(ab)+/</code>	una o più stringhe "ab"	<code>"ab"</code> <code>"abab"</code> <code>"ababab"</code>
<code>/ab+/ </code>	una "a" seguito da una o più "b"	<code>"ab"</code> <code>"abb"</code> <code>"abbb"</code>

Le parentesi tonde **memorizzano** la stringa di testo corrispondente al contenuto delle parentesi:

- la stringa viene memorizzata in una variabile temporanea
- Il contenuto della variabile può essere richiamato con `\<numero>`
- **1** = contenuto della prima coppia di parentesi; **2** = contenuto della seconda coppia di parentesi, ecc.

```
/(le|gli)(il|lo|la)+\1/ \([bcdfghjklmnpqrstvwxyz])+1/
```

Raggruppamento e memoria

RE	Definizione	Esempi di "matching"
<code>/(ab)+\1/</code>	la variabile "\1" corrisponde a qualunque stringa abbia fatto matching con il contenuto delle parentesi	<u>"abab"</u> <u>"abababab"</u> <u>"abababababab"</u>
<code>/(a)+(b)+\1\2/</code>	la variabile "\1" corrisponde a qualsiasi stringa abbia fatto matching con il contenuto della prima coppia di parentesi; la variabile "\2" idem, ma rispetto alla seconda coppia di parentesi	<u>"abab"</u> <u>"abbabb"</u> <u>"aabaab"</u> "abbaab"
<code>/p(.)o p\1o/</code>	la variabile "\1" corrisponde a qualunque stringa abbia fatto matching con il contenuto delle parentesi	<u>"pio pio"</u> <u>"pao pao"</u> <u>"pro pro"</u> , ecc. "pio pao" "pao peo"
<code>/p.o p.o/</code>	la wildcard "." può essere sostituita da qualsiasi carattere	<u>"pio pio"</u> <u>"pio pao"</u> <u>"pro pso"</u> , <u>"pao pio"</u> , ecc.
<code>/(a)(b)+\1\2/</code> <code>/a(b)+a\1/</code>		<u>"abab"</u> , <u>"abbabb"</u> , "aabaab", "abbaab"

Priorità

Ordine	Tipo	Esempi di "matching"
1	parentesi	()
2	moltiplicatori	? + * {m,n} {m,} {n} ?? +? *?
3	Sequenza e ancore	"cane" ^ \$ \b \B
4	alternativa	

Caratteri speciali

- Alcuni caratteri hanno un significato speciale nel linguaggio delle RE
 - `[]?*.()+-/{ }`
- Se questi caratteri fanno parte del pattern di testo da cercare, devono comparire in una RE con davanti il carattere `\` (**carattere di escape**)
 - `/[a/` ERRORE! '[' è interpretato come classe di carattere e manca la parentesi ']'
 - `/\[a/` la stringa "[a"
 - `/a./` qualsiasi stringa di due caratteri che inizia con 'a' "ab" "au" as" "a1" "a?" ...
 - `/a\./` la stringa di testo "a."
 - `/cane?/` le stringhe "cane" e "can"
 - `/cane\?/` la stringa "cane?"

Ancore

- Le **ancore** sono caratteri speciali che specificano dove deve comparire il pattern di testo da cercare
 - `/^<pattern>/` il <pattern> deve comparire all'inizio di una linea
 - `<pattern>$/` il <pattern> deve comparire alla fine di una linea

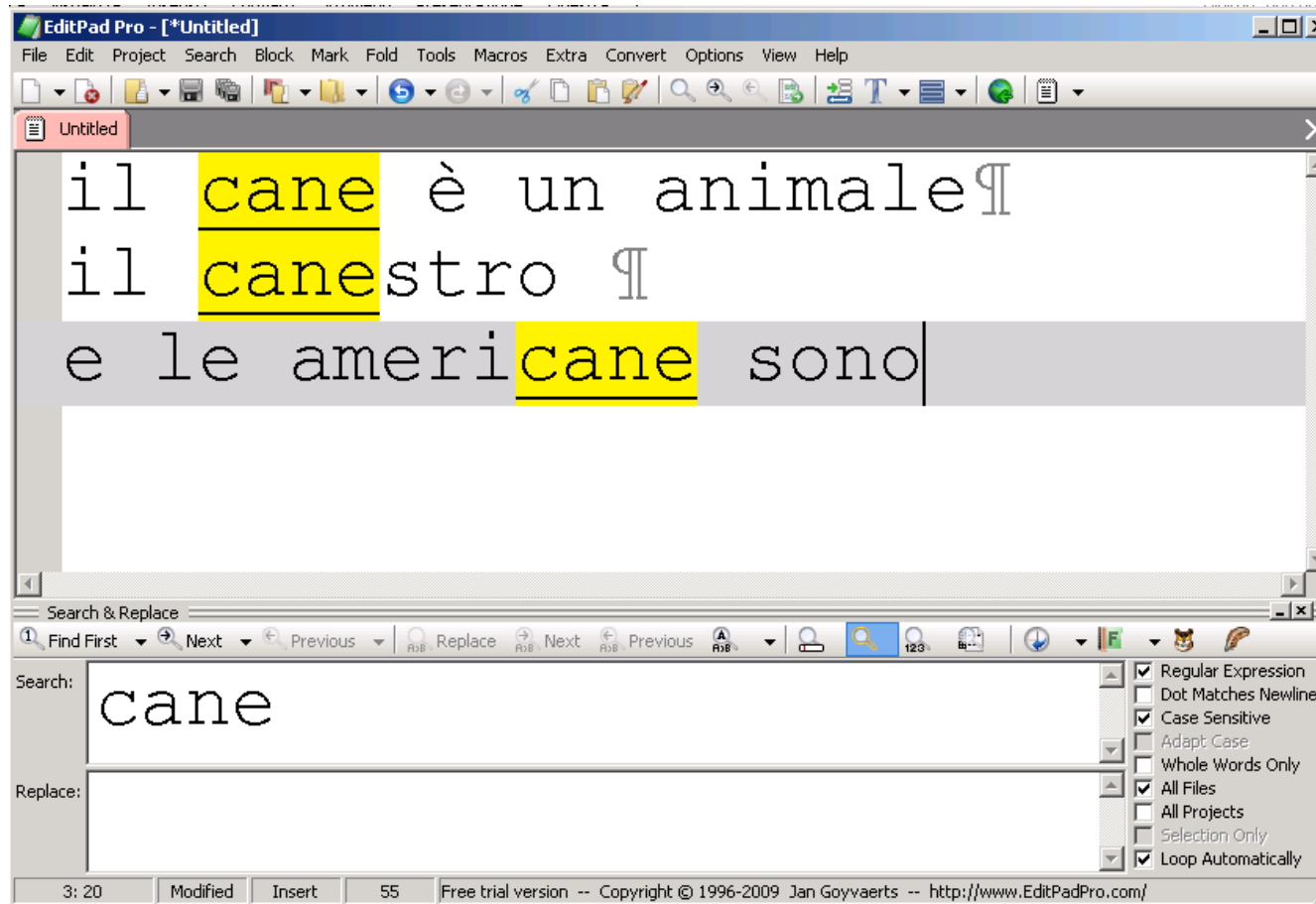
RE	Definizione	Esempi di "matching"
<code>/cane\$/</code>	la stringa 'cane' quando compare alla fine di una linea	<code>"...cane"</code> "il cane di Mario"
<code>/^La/</code>	la stringa 'La' quando compare all'inizio di una linea	<code>"La macchina era guasta"</code> "il treno per La Spezia"
<code>/^La Spezia\$/</code>	una riga che contiene solo la stringa "La Spezia"	<code>"La Spezia"</code> "...a La Spezia per lavoro ..."

Ancore

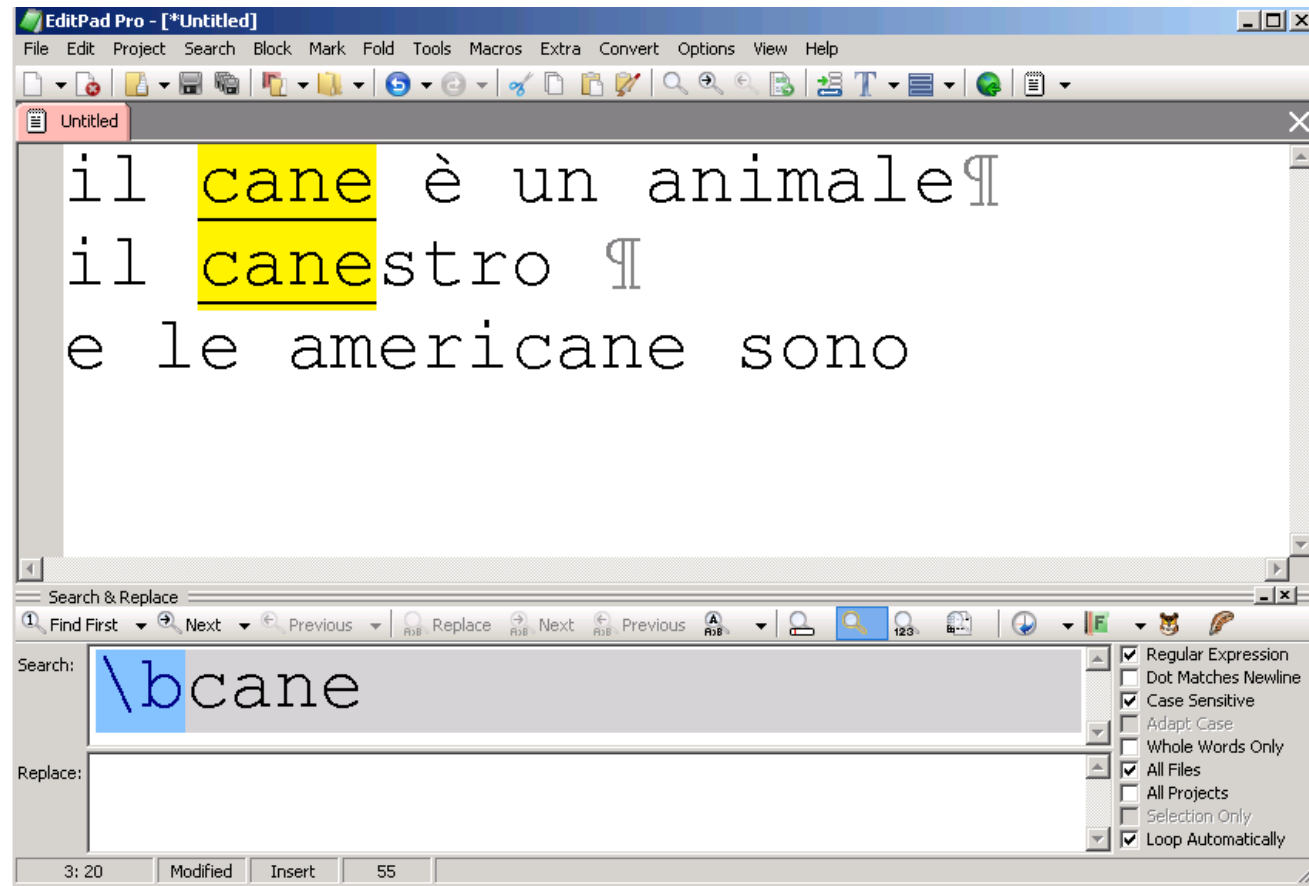
- “\b” è un ancora che indica il confine di una parola (“\B” indica ogni punto non confine di parola)
 - Il confine di una parola è un punto che ha da una parte un carattere di classe \w e dall'altra o un carattere di classe \W o l'inizio (fine) riga
 - **ATTENZIONE:** I caratteri accentati (à, è, è, ì, ò, ù) fanno parte della classe \W (così come lo spazio e gli altri segni di punteggiatura)

RE	Definizione	Esempi di “matching”
<code>/\bcane\b/</code>	la stringa 'cane' deve avere a destra e a sinistra un confine di parola	“il <u>cane</u> è ...” “il canestro” “le americane sono”
<code>/\Bcane\b/</code>	la stringa 'cane' deve avere a destra (ma non a sinistra) un confine di parola	“il cane è ...” “il canestro” “le ameri <u>cane</u> sono”
<code>/\bè\b/</code>	Trovare 'è' come copula	“Pinocchio <u>è</u> stanco”
<code>/\Bè\B/</code>		“Pinocchio è stanco”

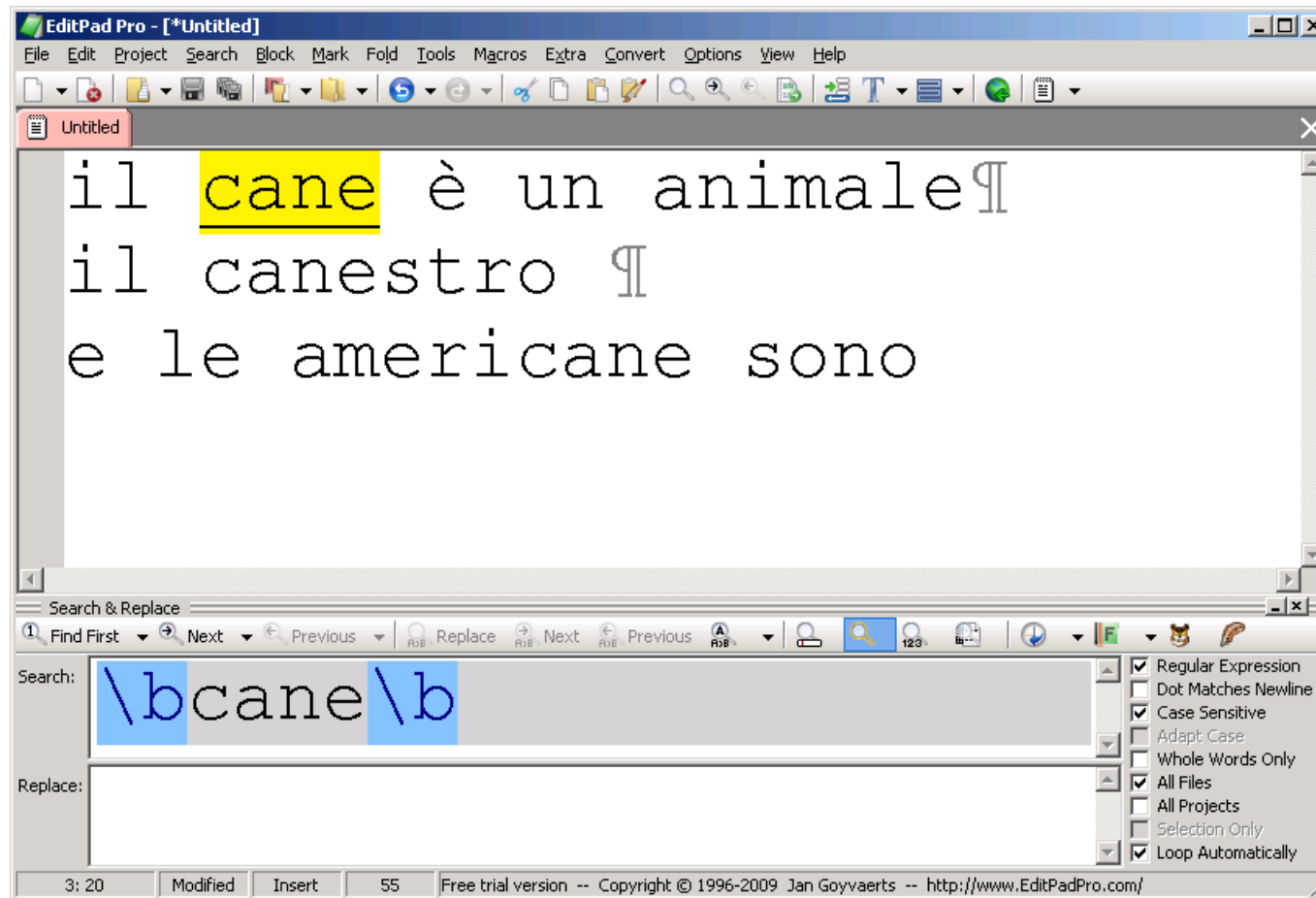
Esempio con EditPad Pro



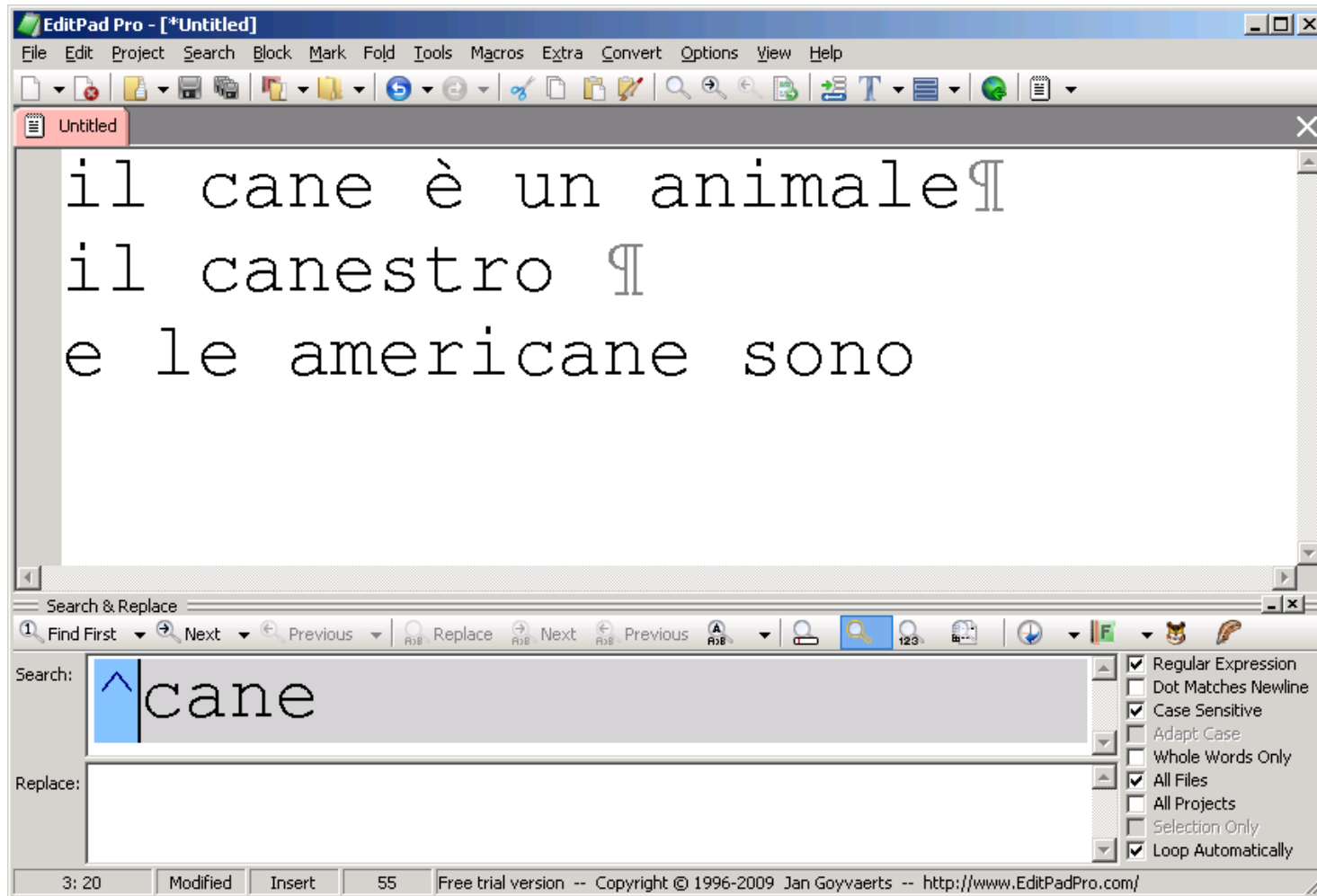
Esempio con EditPad Pro



Esempio con EditPad Pro

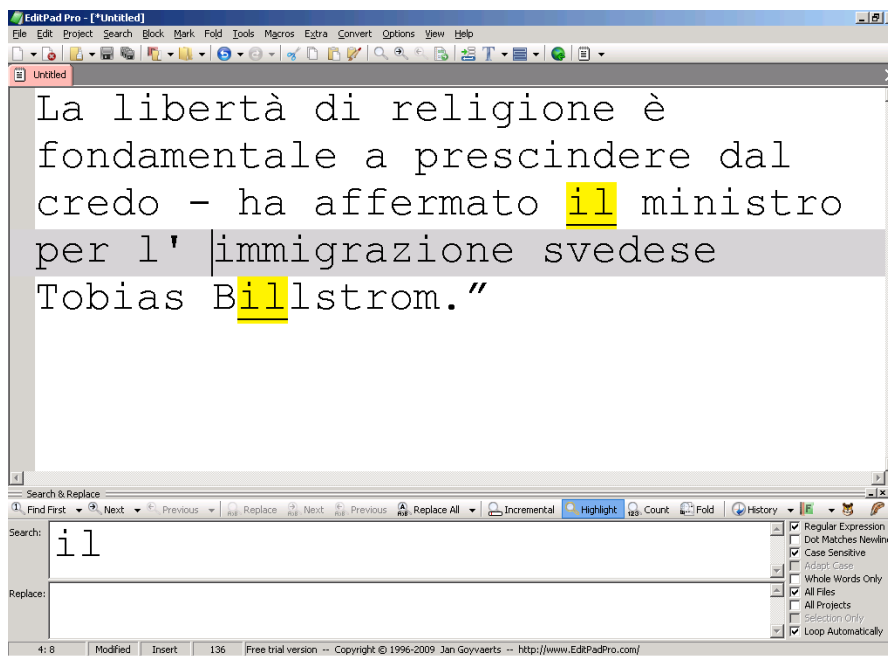


Esempio con EditPad Pro

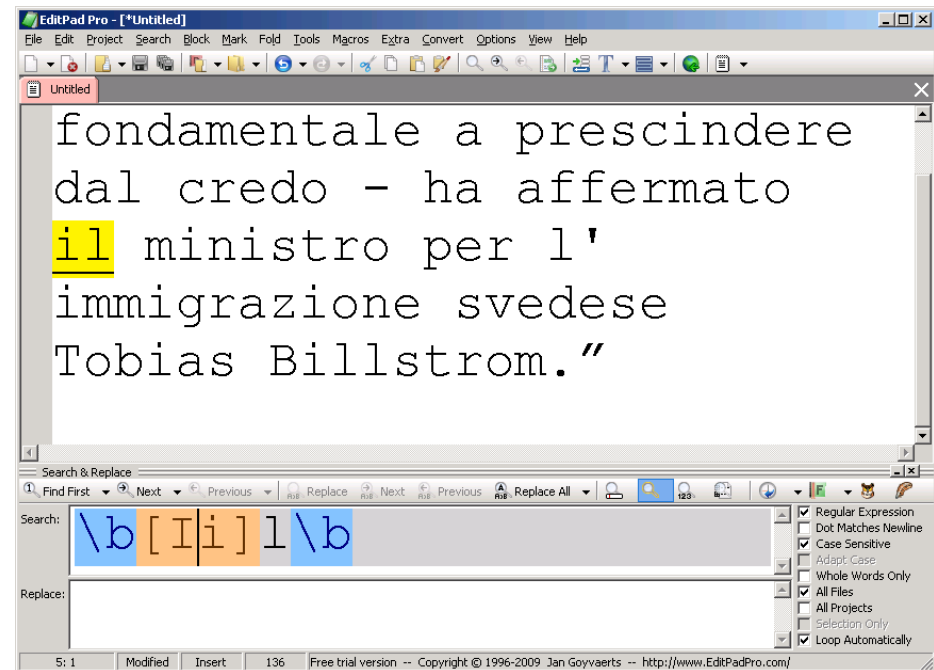


Esempi con EditPad Pro

- Scrivere una espressione regolare per trovare l'articolo "il":



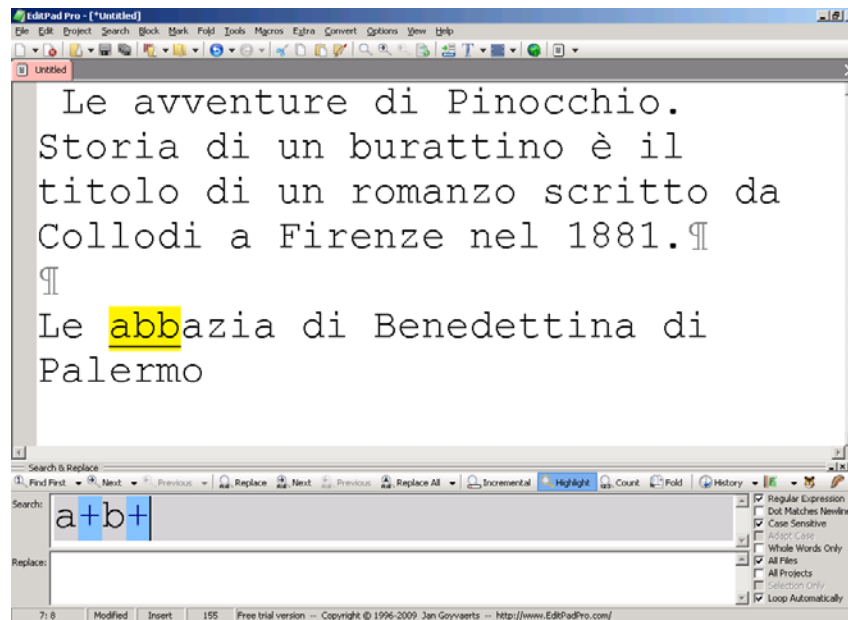
Errata



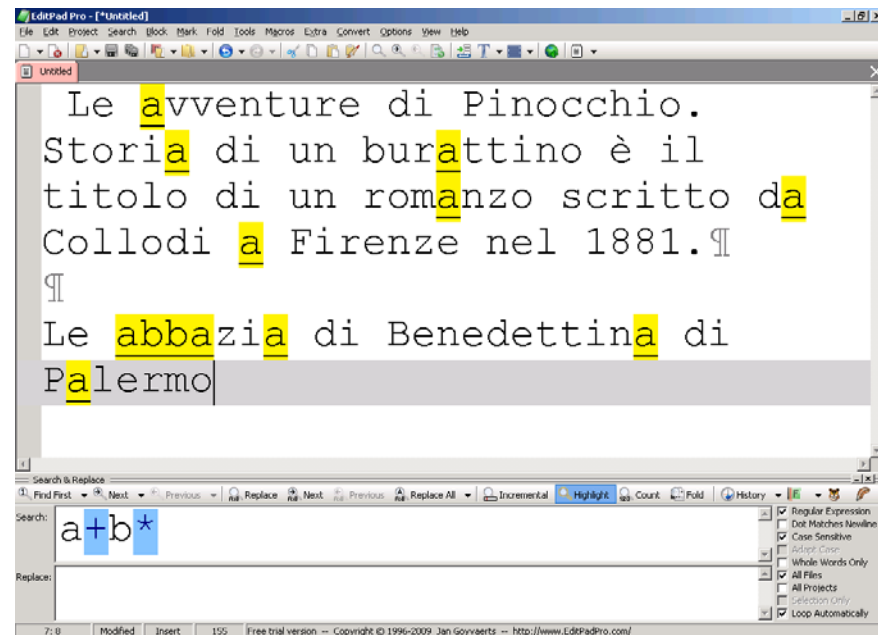
Corretta

Esempi

- Scrivere una espressione regolare che trovi almeno una 'a' seguita da un qualunque numero di 'b'



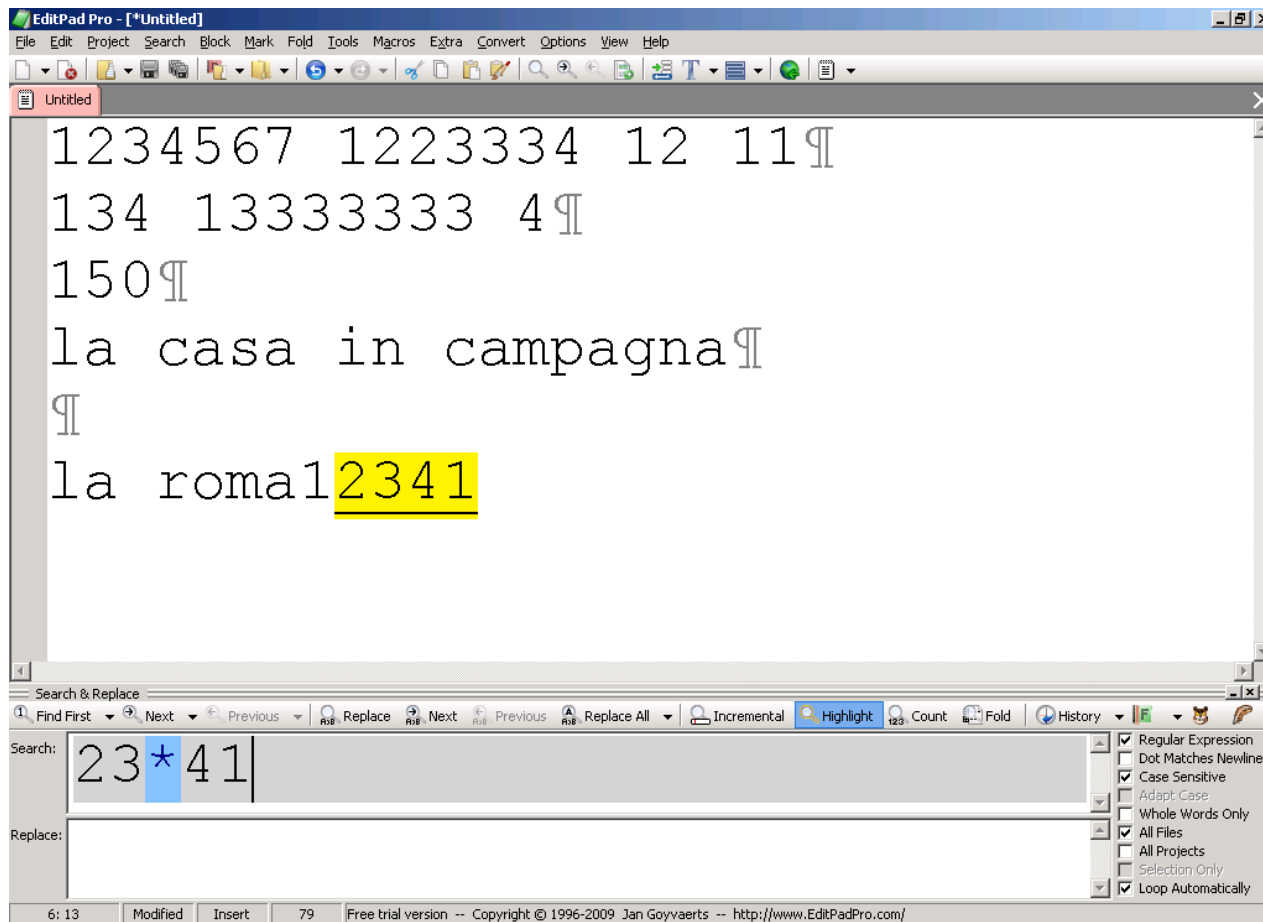
Errata



Corretta

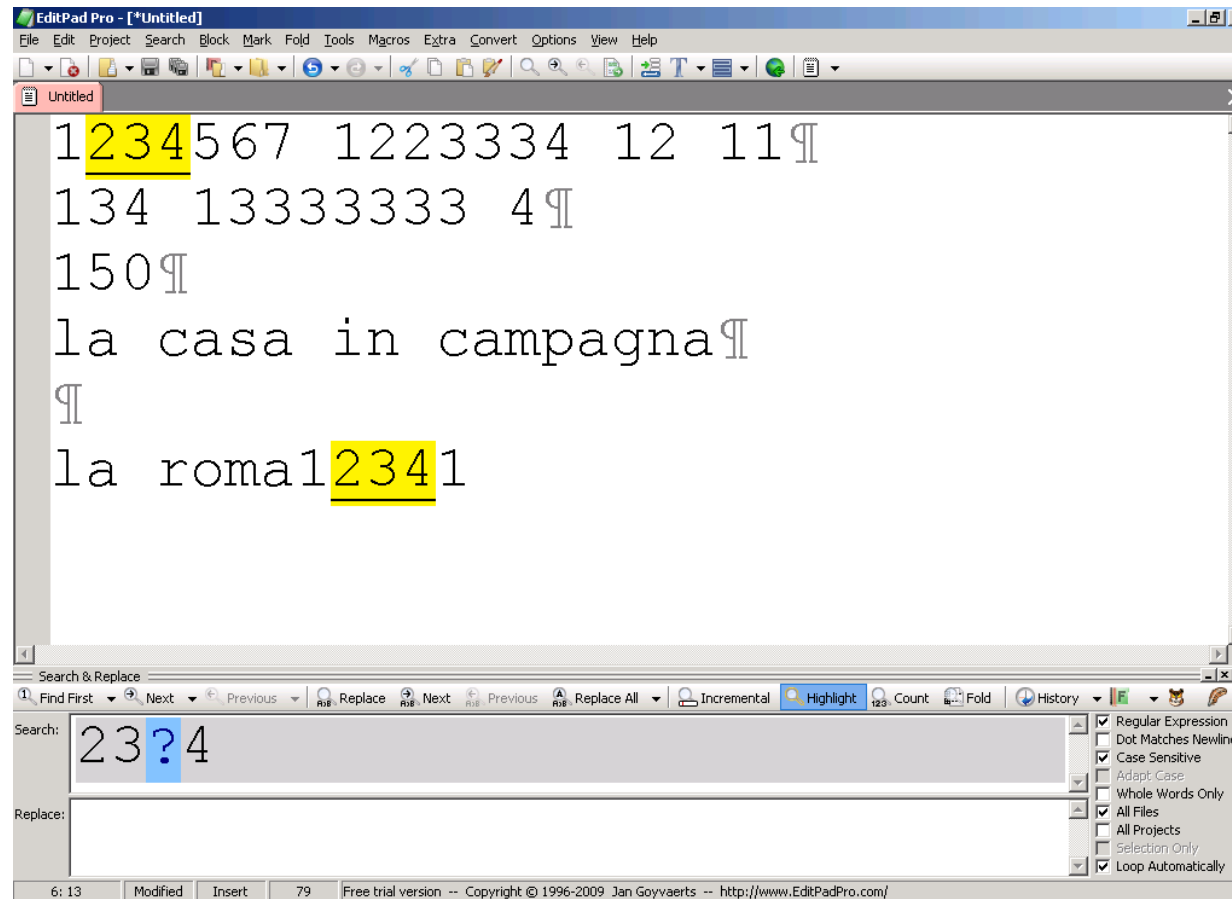
Esempi con EditPad Pro

- Moltiplicatori di carattere



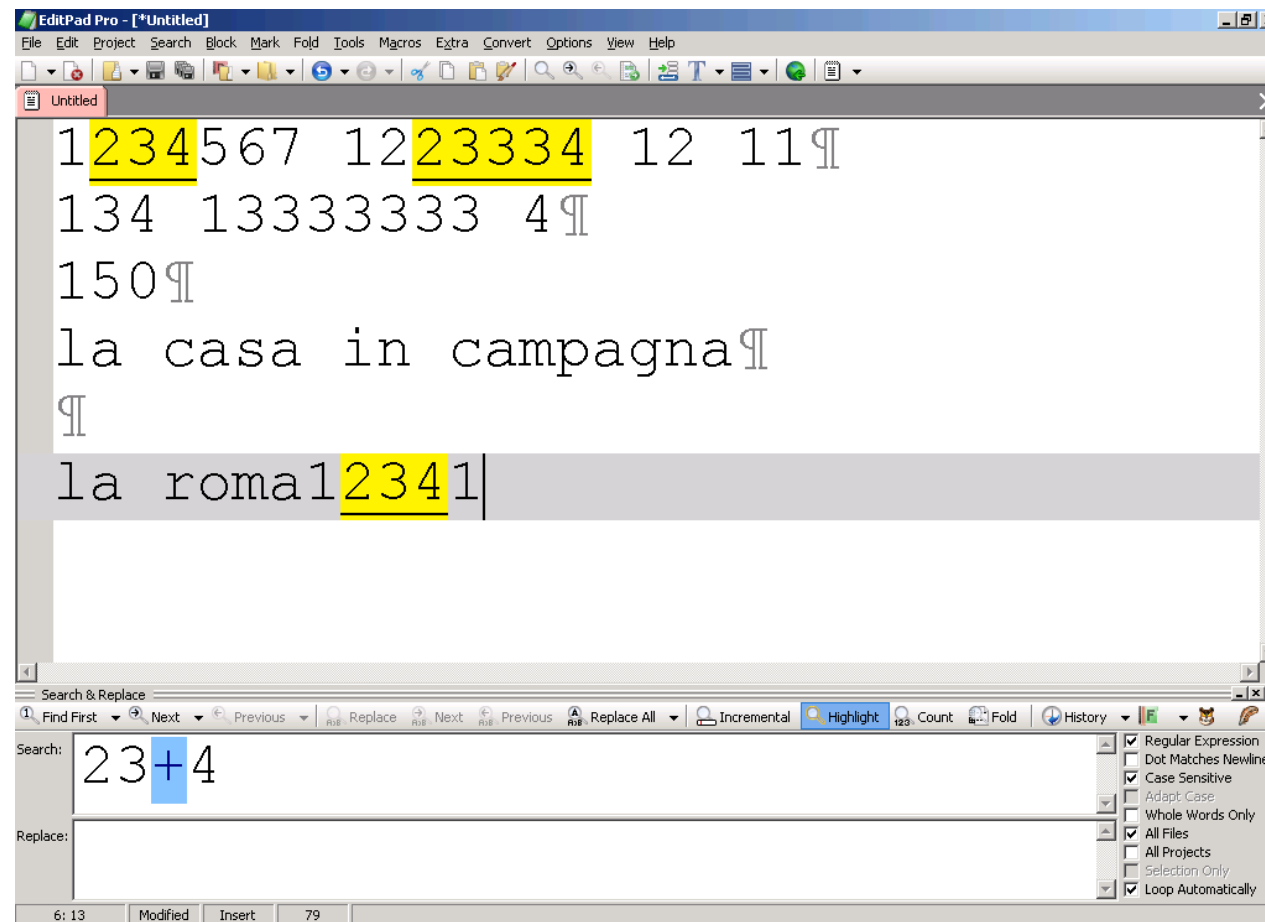
Esempi con EdiPad Pro

- Moltiplicatori di carattere



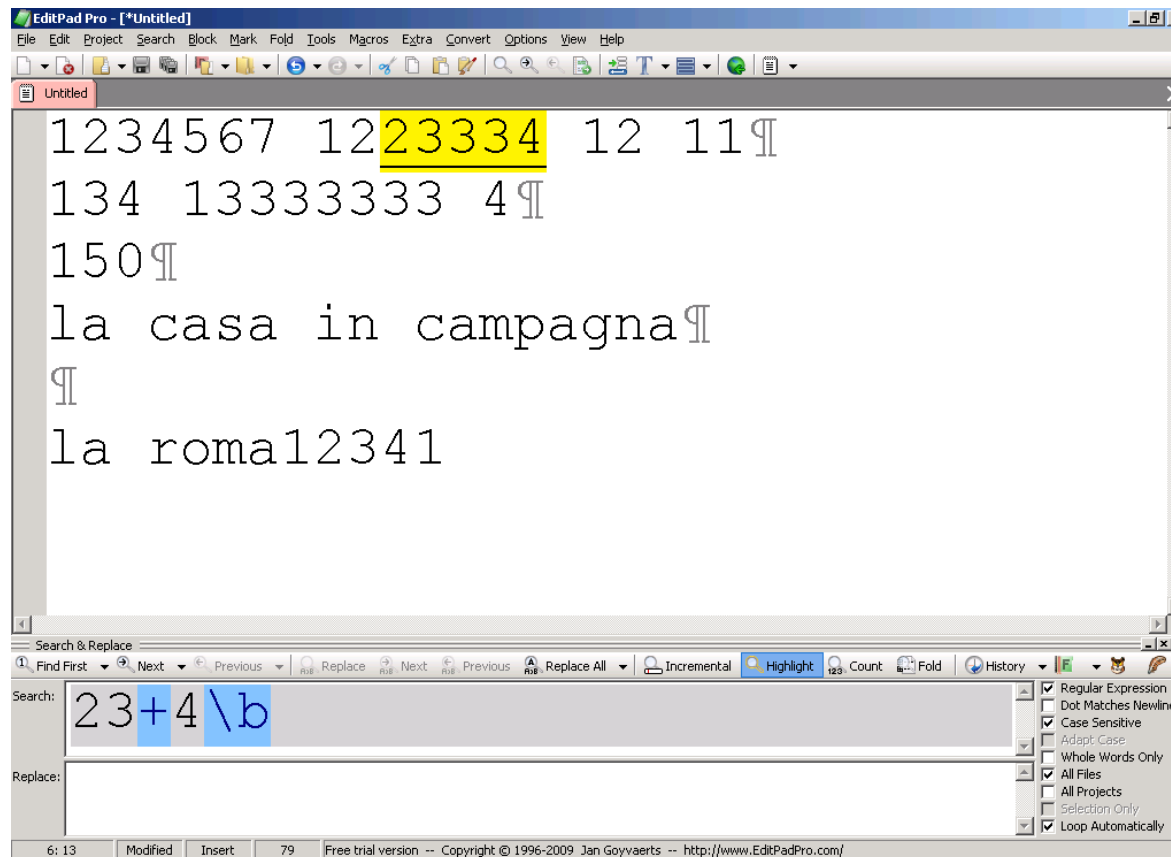
Esempi con EdiPad Pro

- Moltiplicatori di carattere



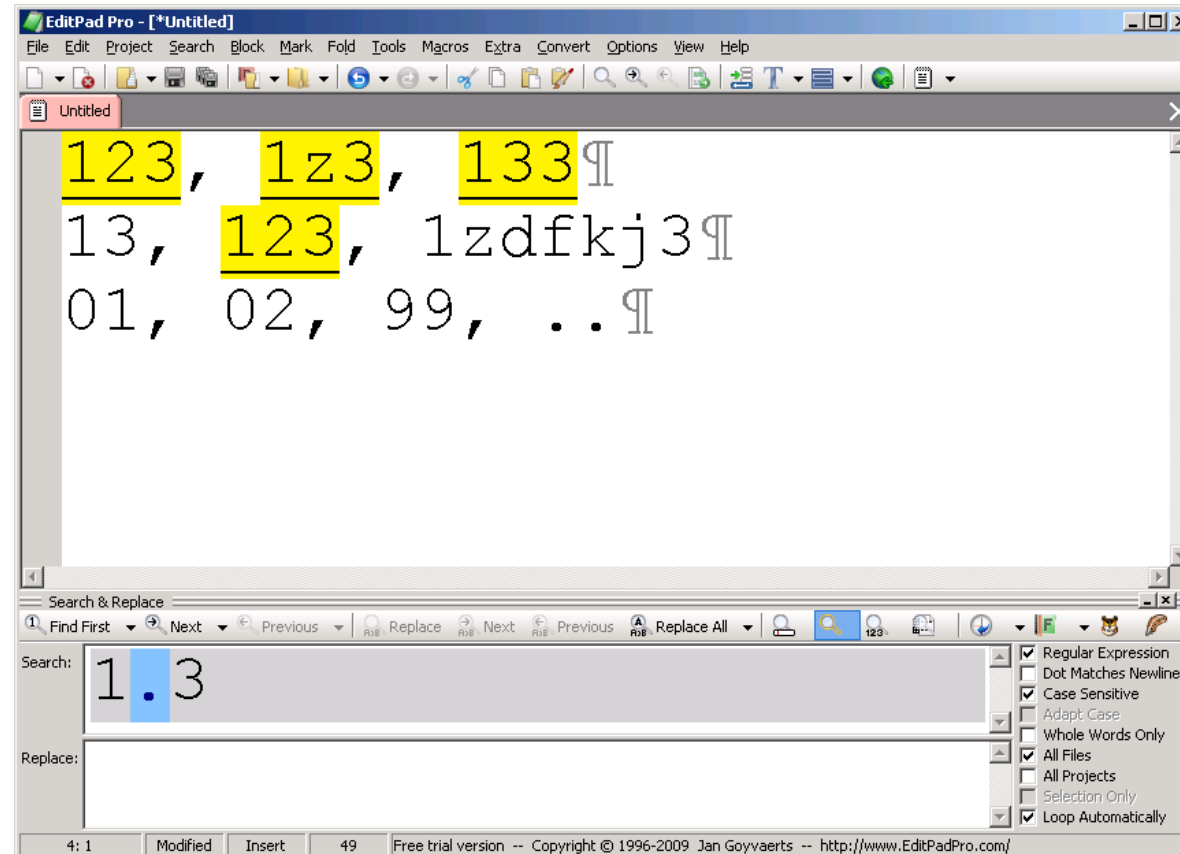
Esempi con EdiPad Pro

- Moltiplicatori di carattere



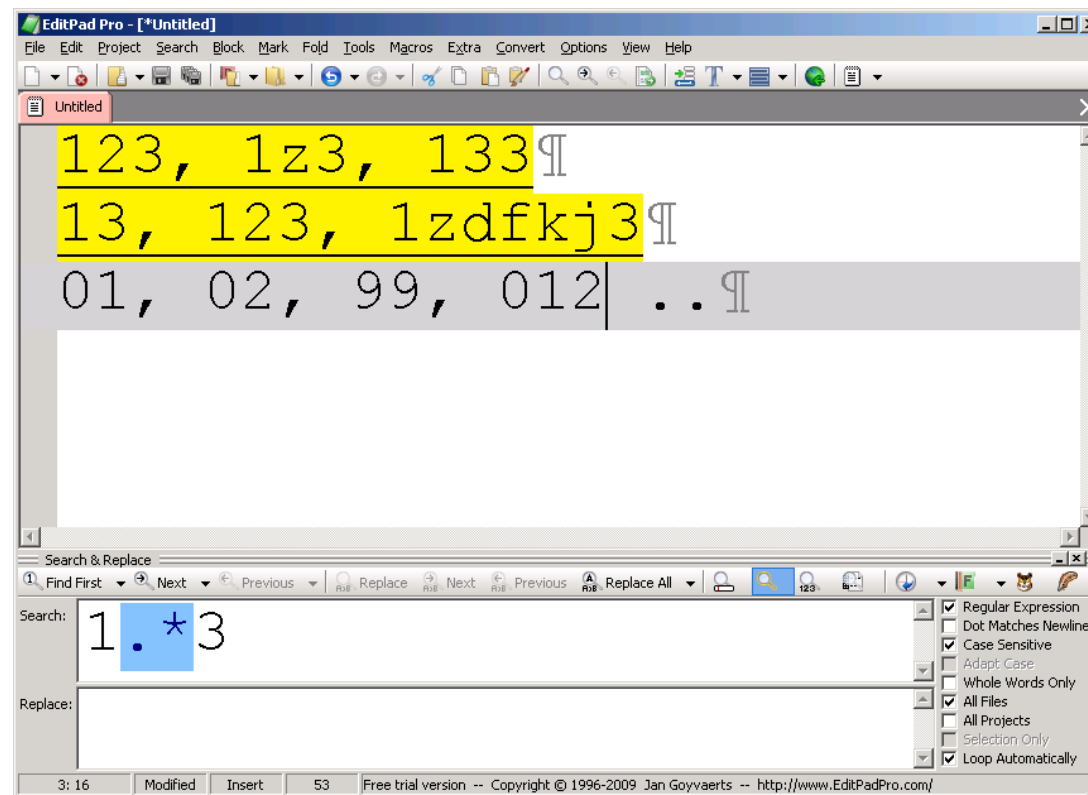
Esempi con EditPad Pro

- ‘.’ qualunque carattere eccetto ritorno a capo



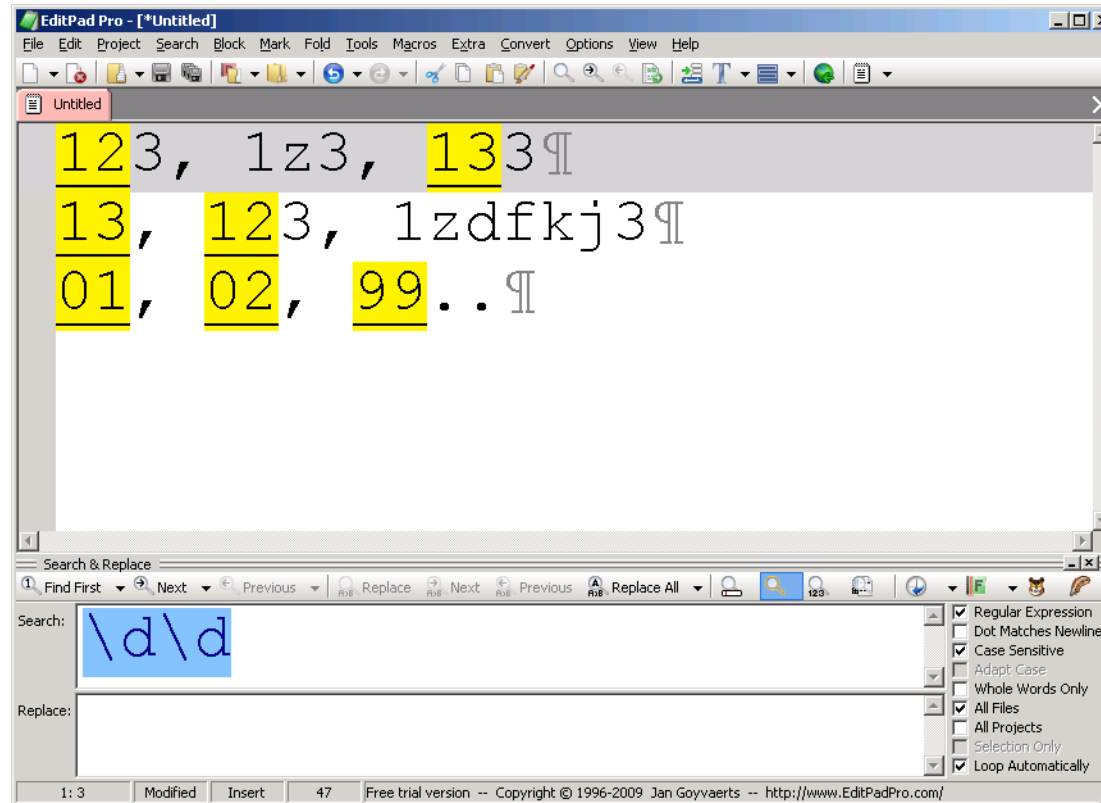
Esempi con EditPad Pro

- ‘.’ qualunque carattere eccetto ritorno a capo



Esempi con EditPad Pro

- due interi tra 0 e 9



Esempi

- Classi di caratteri

RE	Matches	Does Not Match
/[^ab]/	c, d, z	ab
^[1-9][0-9]*\$	Qualunque numero positivo	Zero, numeri negativi o decimali
[0-9]*[.]?[0-9]+	0.1 .8 1 1.2 100,000	12.

Esercizi

- Formalizzare con le espressioni regolari i patterns per trovare le seguenti stringhe
 - “tutte le vocali minuscole o maiuscole”

Esercizi

- Formalizzare con le espressioni regolari i patterns per trovare le seguenti stringhe
 - “tutte le vocali minuscole o maiuscole”

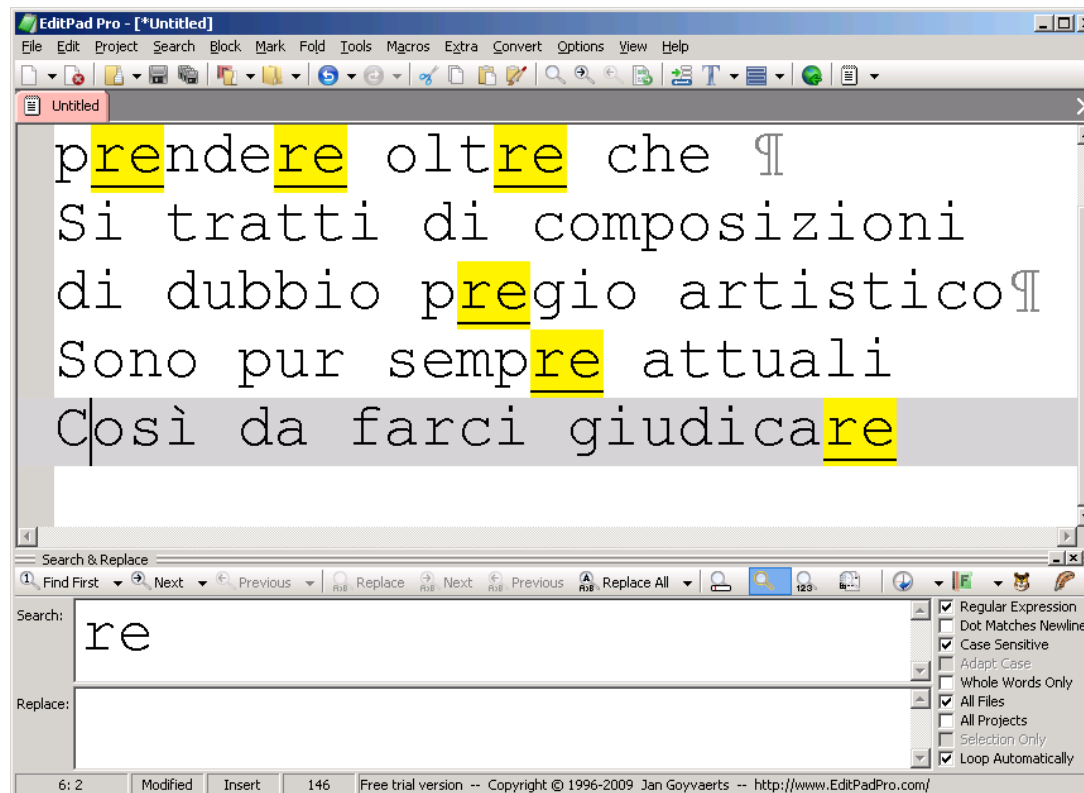
Sol.: `/[AaEeIiOoUu]/`

Esercizi EditPad Pro

- Formalizzare con le espressioni regolari i patterns per trovare le seguenti stringhe
 - 're'

Esercizi EditPad Pro

- Formalizzare con le espressioni regolari i patterns per trovare le seguenti stringhe
 - 're'



Esercizi EditPad Pro

- Formalizzare con le espressioni regolari i patterns per trovare le seguenti stringhe
 - Che terminano in 're'

Esercizi EditPad Pro

- Formalizzare con le espressioni regolari i patterns per trovare le seguenti stringhe
 - Che terminano in 're'
 - /re/
 - /re\b/
 - \bre\b/
 - \Are/
 - /re\$/

 - /^[^re]/
 - /[re]

