



Python per la Manipolazione dei Testi

Parte 2

Cezar Sas

University of Groningen

c.a.sas@rug.nl

UNICODE

(e cenni sulle sue follie)



Categorie Unicode

- Caratteri Unicode (minuscolo/maiuscolo) - **LI** **Lu**
- Operatori Matematici (= - + *)- **Sm**
- Caratteri Valuta (\$ € £)- **Sc**
-

Usabili con:

```
pip install regex  
import regex
```



Categorie Unicode

- `[\p{L} -- QW]` – Tutti i caratteri senza QW
- `[\p{N} -- [0 - 9]]` – Tutti i caratteri numerici senza [0-9]
- . . .

Ottimizzazione di Espressioni Regolari *(Tips & Tricks)*

Backtracking

Comportamento di una regex quando fallisce. Ovvero arretra (backtracking) per provare a valutare possibili alternative.

Esempio:

`r".+b" → acaaacbad`

`acaaacbad` - `.+`

`acaaacbad` - `.+`

...

`acaaacbad` - `.+`

`acaaacbad` - `.+`

BACKTRACKING

`acaaacbad` - `.+b`

BACKTRACKING

`acaaacbad` - `.+b`

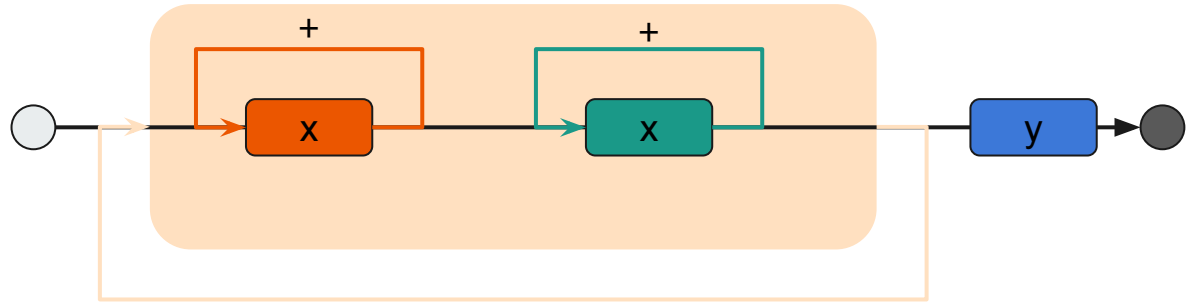
BACKTRACKING

`acaaacbad` - `.+b`

Backtracking Catastrophico

$r''(x+x+)+y''$

xxxxxxxxxy

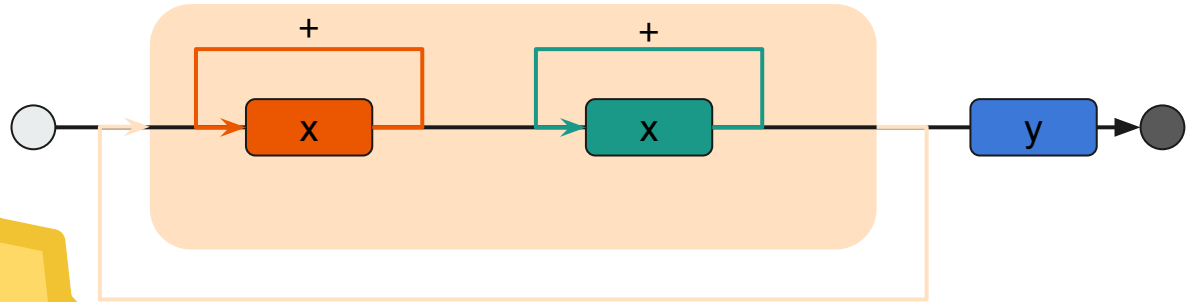


Backtracking Catastrophico

$r''(x+x+)+y''$

xxxxxxxxxy

2558
passi





Backtracking Catastrofico

Problema:

Quantificatori nested che hanno multiple opzioni di matching

Soluzione:

- Semplificare la regex
- *Quantificatori possessivi e gruppi atomici* (disponibili nel modulo `regex`)



Tutto Opzionale

Esempio:

```
[-+]?[0-9]*\.[0-9]*
```

Matcha anche il carattere vuoto (ϵ)

Soluzione:

```
[-+]?([0-9]*\.[0-9]+|[0-9]+)
```

Guida in dettaglio

<https://www.regular-expressions.info/tutorial.html>

LIVE!  REC

CODING

(Soluzione Esercizi)

—