Expresiones regulares

Fernando Carranza fernandocarranza86@gmail.com

Miércoles 20/03/2024

Estructura

En esta charla vamos a tratar de cumplir los siguientes objetivos:

- Hacer un acercamiento al lugar de las expresiones regulares en el marco de la teoría de los lenguajes
- Presentar las convenciones notacionales de las expresiones regulares
- Revisar las funcionalidades de la librería de Python re, para la manipulación de textos mediante expresiones regulares

Teorema de Borges

Teorema de Borges

Un mono tipeando teclas al azar en una máquina de escribir tipeará las obras completas de Shakespeare si le damos suficiente tiempo (Pablo Groisman. Te regalo un teorema. TantaAgua)

La Biblioteca Total

Todo estará en sus ciegos volúmenes. Todo: la historia minuciosa del porvenir, Los egipcios de Esquilo, el número preciso de veces que las aguas de Ganges han reflejado el vuelo de un halcón, el secreto y verdadero nombre de Roma, (...) Todo, pero por una línea razonable o una justa noticia habrá millones de insensatas cacofonías, de fárragos verbales y de incoherencias. Todo, pero las generaciones de los hombres pueden pasar sin que los anaqueles vertiginosos -los anaqueles que obliteran el día y en los que habita el caos- les hayan otorgado una página tolerable.

(Jorge Luis Borges. La Biblioteca total.)

Corpus

- En términos formales, el corpus que esecribió el mono es en esencia una larga cadena de símbolos.
- Ese mono es capaz de escribir cosas muy complejas, pero es incapaz de distinguir esas cosas de otras que son fárrragos verbales.

• Supongamos que queremos extraer ciertos datos de ese corpus. Nuestro deseo puede representarse en términos de un *problema*.

- Supongamos que queremos extraer ciertos datos de ese corpus.
 Nuestro deseo puede representarse en términos de un problema.
- A cada problema, un algoritmo.

Algunos problemas para nuestro corpus

- ¿Cuáles son los guarismos escritos en números arábigos?
- ¿Cuáles son aquellas cadenas que expresan operaciones aritméticas?
- ¿Cuáles son aquellas operaciones aritméticas bien hechas y cuáles no?
- ¿Cuáles son las obras de Shakespeare?
- ¿Cuáles son los textos escritos en español?

Algunos problemas para nuestro corpus

- ¿Cuáles son los guarismos escritos en números arábigos?
- ¿Cuáles son aquellas cadenas que expresan operaciones aritméticas?
- ¿Cuáles son aquellas operaciones aritméticas bien hechas y cuáles no?
- ¿Cuáles son las obras de Shakespeare?
- ¿Cuáles son los textos escritos en español?

No todos los problemas acarrean la misma *complejidad*. La complejidad del algoritmo para solucionarlo dependerá de la complejidad del problema.

Un chiste nerd



Un poco de terminología

• Denominamos alfabeto al conjunto no vacío de símbolos discretos. En el caso de los textos digitales, estos serán caracteres pertenecientes a alguna fuente de codificación. El alfabeto se designa convencionalmente con la letra Σ

Un poco de terminología

- Denominamos alfabeto al conjunto no vacío de símbolos discretos. En el caso de los textos digitales, estos serán caracteres pertenecientes a alguna fuente de codificación. El alfabeto se designa convencionalmente con la letra Σ
- La concatenación de símbolos (iguales o diferentes) de un alfabeto Σ se conoce con el nombre de cadena.



 Σ^* equivale al conjunto de todas las cadenas posibles que se pueden obtener a partir de un alfabeto $\Sigma.$

Lenguaje

 Denominamos lenguaje a un conjunto de cadenas que está incluido en Σ*. Definir un lenguaje por extensión consiste en nombrar uno a uno todos sus miembros. Esto, por supuesto, es un método de definición de conjuntos totalmente inviable para lenguajes formados por infinitos elementos. En su lugar, hay que recurrir a un tipo de definición intensional.

- Todo problema (e.g., ¿cuáles cadenas pertenecen al español?, ¿cuáles son las obras de Shakespeare?, etc.) se puede representar en términos de un lenguaje (e.g. el conjunto de cadenas del español, el conjunto de las obras de Shakespeare, etc.) y viceversa.
- Así como hay problemas de diversa complejidad, los respectivos lenguajes también se clasifican según su complejidad.

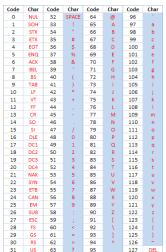
 Las expresiones regulares son un sistema de notación que permite caracterizar la familia de lenguajes más sencillos: los llamados lenguajes regulares.

- Las expresiones regulares son un sistema de notación que permite caracterizar la familia de lenguajes más sencillos: los llamados lenguajes regulares.
- Esa sencillez formal acarrea también un bajo costo de procesamiento, lo que vuelve a las expresiones regulares muy eficientes computacionalmente.

- Las expresiones regulares son un sistema de notación que permite caracterizar la familia de lenguajes más sencillos: los llamados lenguajes regulares.
- Esa sencillez formal acarrea también un bajo costo de procesamiento, lo que vuelve a las expresiones regulares muy eficientes computacionalmente.
- Las expresiones regulares no sirven, entonces, para caracterizar conjuntos de cadenas más complejos, como el conjunto de las oraciones gramaticales del español o el conjunto de operaciones aritméticas bien resueltas. Pero permiten identificar lenguajes más sencillos, como el conjunto de las cadenas de caracteres alfanuméricos separados por puntos, el conjunto de las cadenas que tienen el formato de direcciones de correo electrónico, etc.

- Las expresiones regulares son un sistema de notación que permite caracterizar la familia de lenguajes más sencillos: los llamados lenguajes regulares.
- Esa sencillez formal acarrea también un bajo costo de procesamiento, lo que vuelve a las expresiones regulares muy eficientes computacionalmente.
- Las expresiones regulares no sirven, entonces, para caracterizar conjuntos de cadenas más complejos, como el conjunto de las oraciones gramaticales del español o el conjunto de operaciones aritméticas bien resueltas. Pero permiten identificar lenguajes más sencillos, como el conjunto de las cadenas de caracteres alfanuméricos separados por puntos, el conjunto de las cadenas que tienen el formato de direcciones de correo electrónico, etc.
- Aplicado como problema, una expresión regular nos permite identificar cadenas de determinada naturaleza en el marco de un corpus, como el que elaboró el mono.

Primeramente, cualquier secuencia de caracteres del alfabeto $\Sigma = \mathsf{ASCII}$, escapando aquellos caracteres que oficien de operadores, es una expresión regular r que matcheará con esa secuencia de caracteres.



En la bibliografía teórica se utiliza la función interpretación $[\![\alpha]\!]$ para evaluar el significado de la expresión regular α (ver Wintner 2010). Por ejemplo,

- $[abc] = \{w | w \text{ es una cadena formada por la concatenación de a, b y } c\} = abc$
- y así sucesivamente.

Paréntesis y estrella de Kleene

- ()
 Para delimitar bloques de cadenas
- cero o más ocurrencias del caracter o bloque anterior

 Por ejemplo, [a*] = {w| w es una cadena formada por la concatenación de cero o más ocurrencias de a}

```
[[ab*]]
```

+y?

- una o más ocurrencias del caracter o bloque anterior

 Por ejemplo, [a+] = {w| w es una cadena formada por la concatenación de una o más ocurrencias de a}
- !
 una o ninguna ocurrencia del caracter o bloque anterior
 Por ejemplo, [ab?] = {a, ab}

$[\![ab+]\!]$	
$[\![(ab)+]\!]$	
[(a?b)+]	

- - cualquier caracter excepto cambio de línea Por ejemplo, $[a.] = \{aa, ab, ac, ad...\}$
- {n}
 n cantidad de ocurrencias del caracter o bloque anterior
 Por ejemplo, [a{3}] = {aaa})
- {n,m}
 entre n y m cantidad de ocurrencias del caracter o bloque anterior (cuando se usa para matchear, matchea con el menor número posible)
 Por ejemplo, [a{3,6}] = {aaa, aaaa, aaaaa, aaaaaa})

[Achi{2}s] [(Bla){3}] [(Bla){2,4}]

Disyunción y rango

- $\bullet \ [] = {\sf disyunci\'on\ entre\ lo\ que\ aparezca\ adentro\ (\llbracket [{\sf abcd}] \rrbracket = \{{\sf a,b,c,d}\}) }$
- [x-y] = rango de caracteres que aparecen entre x e siguiendo el orden de ASCII ($[[a-d]] = \{a,b,c,d\}$)
- $[^{\land}x]$ = expresiones que excluyan x.

El rango mayor que se puede obtener es $[-\sim]$

Abreviaturas

- $\sl s =$ espacio en blanco
- \bullet \S = cualquier caracter que no sea un espacio en blanco
- \d = cualquier dígito.
- $\backslash D$ = cualquier símbolo que no sea un dígito
- \w = cualquier caracter alfanumérico

Establecer las equivalencias en términos de rangos

Abreviatura	rango
\d	
\w	

Expresiones target	Expresión regular
Edad (en arábigos)	
Estructura de sílaba	
Perífrasis terminativas	

Páginas para probar expresiones regulares:

- https://regex101.com/
- https://regexr.com/

Bibliografía I

- Hopcroft, J. E., Motwani, R., y Ullman, J. D. (2006). *Automata theory, languages, and computation*. Addison-Wesley, Boston, Massachusetts.
- Jurafsky, D. y Martin, J. H. (2024). Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition. Draft disponible.
- Wintner, S. (2010). Formal language theory. En Clark, A., Fox, C., y Lappin, S., editores, *The Handbook of Computational Linguistics and Natural Language Processing*, pp. 11–42. Willey Blakwell, Massachusetts.