

2 - Tareas Básicas PLN: NLTK y Spacy

Tecnologías Emergentes en la Sociedad de
la Información

NLTK

- Es una librería de Python para realizar tareas de PLN.
 - División de oraciones y tokenización.
 - Análisis Morfosintáctico.
 - Lematización y stemming.
 - Reconocimiento de Entidades
 - Expresiones regulares.

Spacy

- Es una librería de Python para realizar tareas de PLN.
 - Tokenización
 - Análisis Morfosintáctico
 - Reconocimiento de Entidades
 - Análisis de Dependencias
 - Similitud Semántica
 - Word Embeddings
- <https://nlpforhackers.io/complete-guide-to-spacy/>

2 - Tareas Básicas PLN: Expresiones Regulares

Tecnologías Emergentes en la Sociedad de
la Información

Expresiones Regulares

- Son patrones que se utilizan para buscar determinadas combinaciones de caracteres en un texto.
- Ej: Encontrar todas las posibles variaciones de la palabra niño (niño, niña, niños, niñas, Niños, Niñas) en un texto.

Expresiones Regulares

- [] el patrón mapea cualquiera de los caracteres contenidos entre corchetes:

Patrón	Qué detecta
[012345]	Cualquier dígito entre 0 y 5
[Aa]	La letra A o a
[nN]iñ[oa]	niño,niña,Niño,Niña

<https://www.regexpal.com/>

Expresiones Regulares

- Caracteres dentro de [] significa que el patrón detectará cualquier texto que contenga algunos de los caracteres entre los corchetes:

Patrón	Qué detecta
[012345]	Cualquier dígito entre 0 y 5
[Aa]	La letra A o a
[nN]iñ[oa]	niño,niña,Niño,Niña

<https://www.regexpal.com/>

Expresiones Regulares

- Rangos [A-Z]

<https://www.regexpal.com/>

Patrón	Qué detecta
[0-9]	Cualquier dígito entre 0 y 9
[A-Z]	Cualquier letra mayúscula
[a-z]	Cualquier letra minúscula
[A-Za-z]	Cualquier letra mayúscula o minúscula

Expresiones Regulares

- | : significa o

Patrón	Qué detecta
- /	- o /
a b c = [abc]	letra a b o c

<https://www.regexpal.com/>

Expresiones Regulares

- | : significa o

Patrón	Qué detecta
- /	- o /
a b c = [abc]	letra a b o c

<https://www.regexpal.com/>

Expresiones Regulares

Patrón	Significado	Ejemplo
?	el carácter anterior es opcional	<i>colou?r -> color, colour</i>
*	0 o más ocurrencias	<i>ja* -> j, ja, jaa, jaa, etc</i>
+	1 o más ocurrencias	<i>ja+ -> j, ja, jaa, jaa, etc</i>
.	cualquier carácter	<i>l.s -> los, las, l-s, ..</i>

<https://www.regexpal.com/>

Expresiones Regulares

- ^ representa el principio de la línea.
- \$ el final de línea
- \ permite representar caracteres especiales como ., *, +

Patrón	Significado
^[0-9]	Identifica las líneas que empiezan con un dígito
;\$	identifica el texto que termina con un ;
\.\$	indentifica el texto que termina con un .

<https://www.regexpal.com/>

Expresiones Regulares

Patrón	Significado
<code>\w</code>	mapea cualquier letra, dígito o underscore
<code>\W</code>	cualquier carácter que no sea <code>\w</code>
<code>\d</code>	cualquier dígito
<code>\s</code>	un espacio en blanco
<code>\t</code>	tabular
<code>\n</code>	salto de línea

<https://www.regexpal.com/>

Ejemplos

- El documento nacional (cédula) de Venezuela tiene el siguiente formato: la letra V (para indicar que es venezolano) o la letra E, un guión, seguido de un número de 8 dígitos.
- Escribe el patrón que permita identificar este tipo de expresiones.

Ejemplo

Regular Expression

JavaScript

flags

```
/[VE]-[\d]{8}/g
```

2 matches

Test String

En Venezuela, al documento nacional de identidad (DNI) se le llama cédula de identidad o simplemente cédula, y generalmente consiste en: una letra (V o E) que indica si se es venezolano o extranjero y un número entre 1 y 99 999 999, separados por un guión, como por ejemplo: **V-15000123** ó **E-82055055**. Sin embargo, el patrón no va a identificar un NIF español como el 08927350-S.

Ejemplo

Escribe la expresión regular para identificar DNI españoles.
Estos pueden ser:

- NIF (Número de Identificación Fiscal) - 8 números y una letra¹
- NIE (Número de Identificación de Extranjeros) - 1 letra², 7 números y 1 letra¹

1 - Una de las siguientes: TRWAGMYFPDXBNJZSQVHLCKE

2 - Una de las siguientes: XYZ

Ejemplo

```
/[\d]{8}-[TRWAGMYFPDXBNJZSQVHLCKE] | [XYZ]-[\d]{7}-[TRWAGMYFPDXBNJZSQVHLCKE]/g
```

2 matches

Test String

V-15000123 ó E-82055055 son ejemplos de cédulas venezolanas. Sin embargo, el patrón no va a identificar un NIF español como el 00000015-P y este un dni de un extranjero X-1234567-S.

```
[\d]{8}-[TRWAGMYFPDXBNJZSQVHLCKE]|[XYZ]-[\d]{7}-[TRWAGMYFPDXBNJZSQVHLCKE]
```

Ejemplo

Regular Expression

Javascript

flags

```
/([XYZ]-[\d]{7}|[\d]{8})-[TRWAGMYFPDXBNJZSQVHLCKE]/g
```

3 matches

Test String

V-15000123 ó E-82055055 son ejemplos de cédulas venezolanas. Sin embargo, el patrón no va a identificar un NIF español como el 00000015-P y este un dni de un extranjero X-1234567-S. Sin embargo, esté no se reconoce, X-00000015-P. Tampoco A-1234567-S

```
([XYZ]-[\d]{7}|[\d]{8})-[TRWAGMYFPDXBNJZSQVHLCKE]
```

Conclusiones

- Expresiones regulares suelen utilizarse como enfoques iniciales para abordar tareas de PLN.
- También pueden ser utilizadas como features (características) dentro de los algoritmos de aprendizaje automático (por ejemplo, si un texto contiene una determinada expresión regular o no).