

# Procesamiento de Lenguaje Natural{

[NLP]

<Preprocesamiento de Texto>

}

```
1  La Agenda de hoy {
2
3      01  Preprocesamiento de texto
4          <Por que se necesita?>
5
6          02  Tokenizar, lematizar, stemmizar
7              <Unificando el formato de nuestro texto>
8
9              03  Todo en una misma bolsa
10                  <Pipelines de preprocesamiento>
11
12  }
13
14
```

# Preprocesamiento de Texto{



**Ivan Rob** @IchBinAivanRob · Ahora

Sonido pirata es el mejor sonidero de todos!! #Pirata #CumbionLoco  
#Mediometro

=

?



**Ivan Rob** @IchBinAivanRob · Ahora

sonido pirata es el mejor Sonidero de todos!!! 🤪 #Pirata #CumbionLoco  
#0.54Yardas

...

}

# Preprocesamiento de Texto{



**Ivan Rob** @IchBinAivanRob · Ahora

Sonido pirata es el mejor sonidero de todos!! #Pirata #CumbionLoco  
#Mediometro

Mayusculas y minusculas

≠

Presencia de emojis



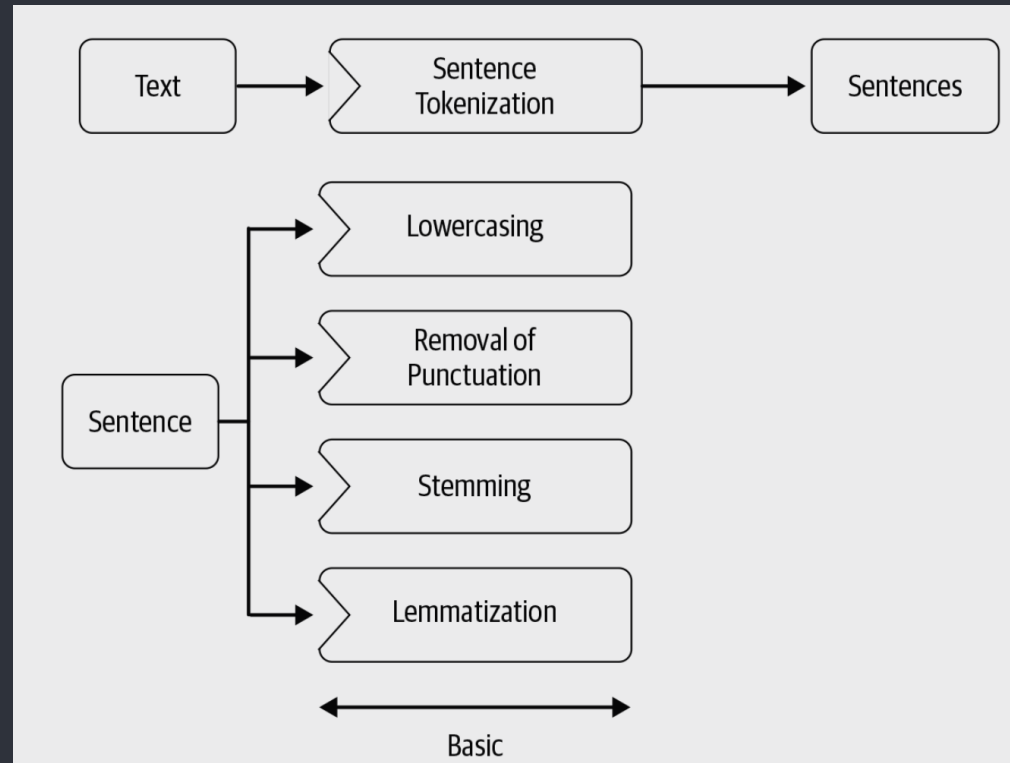
**Ivan Rob** @IchBinAivanRob · Ahora

sonido pirata es el mejor Sonidero de todos!!! 🤪 #Pirata #CumbionLoco  
#0.54Yardas

Diferentes hashtags

# Preprocesamiento de Texto{

```
<Limpiar el texto con el que estaremos trabajando para  
que el análisis/modelado dé los mejores resultados  
posibles. >
```



1  
2  
3  
4  
5  
6  
7 {Lowercasing}  
8  
9  
10  
11  
12  
13  
14

## 1 Lowercasing{

2  
3 < Convierte todos los caracteres de un string a letras  
4 minúsculas>

```
5 nombre = "EL PRINCIPE DEL TACO"  
6 nombre_minus = nombre.lower()  
7 print(nombre_minus)
```

```
8  
9  
10  
11 el principe del taco  
12  
13  
14 }
```

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14

# {Tokenizacion}



# Tokenizacion{

<Dividir una oración en las palabras que la componen. >

```
texto = '''Mi Homero no es comunista. Podrá ser  
mentiroso, puerco, idiota, comunista,  
pero nunca una estrella de porno'''
```

```
texto.split()
```

✓ 0.1s

```
['Mi',  
'Homero',  
'no',  
'es',  
'comunista.',  
'Podrá',  
'ser',  
'mentiroso,',  
'puerco,',  
'idiota,',  
'comunista,',  
'pero',  
'nunca',  
'una',  
'estrella',  
'de',  
'porno']
```

Unigramas



# Tokenizacion{

<Dividir una oración en las palabras que la componen. >

```
import nltk
nltk.download('punkt')
sentence = 'Mi hijo no es comunista podrá ser tonto, estúpido, inútil, comunista pero nunca una estrella porno'
tokens = nltk.word_tokenize(sentence)
print(tokens)
```

✓ 0.0s

```
['Mi', 'hijo', 'no', 'es', 'comunista', 'podrá', 'ser', 'tonto', ',', 'estúpido', ',', 'inútil', ',', 'comunista', 'pero', 'nunca', 'una', 'estrella', 'porno']
```

[nltk\_data] Downloading package punkt to

[nltk\_data] [C:\Users\ivani\AppData\Roaming\nltk\\_data...](C:\Users\ivani\AppData\Roaming\nltk_data...)

[nltk\_data] Package punkt is already up-to-date!

}

Las comas y puntuaciones se toman como su propio token

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14

{Stopwords}

# Stopwords{

< Palabras que no aportan significado o contexto a la oracion>

```
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
stop_words_en = stopwords.words('English')
stop_words_es = stopwords.words('Spanish')
```

```
print(stop_words_en)
```

[ 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', "she's", 'it', 'its', 'they', 'them', 'theirs', 'theirself', 'hers', 'herself', 'us', 'ours' ]

+ Code

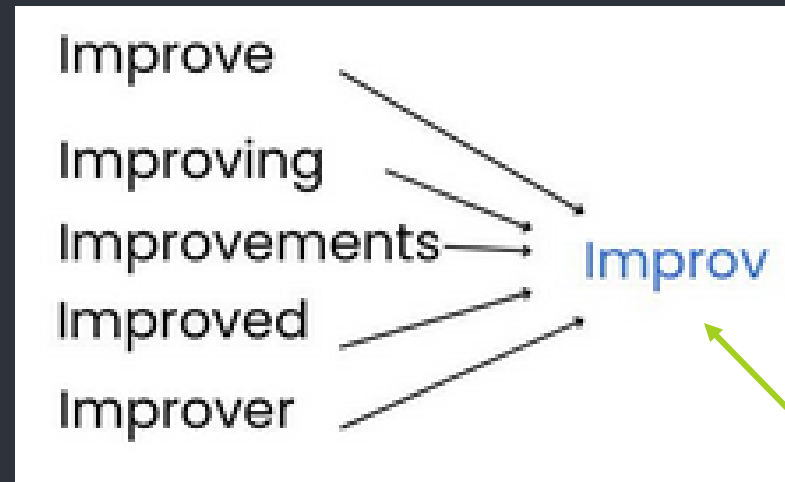
+ Markdown

```
print(stop_words_es)
```

[ 'de', 'la', 'que', 'el', 'en', 'y', 'a', 'los', 'del', 'se', 'las', 'por', 'un', 'para', 'con', 'no', 'una', 'su', 'al', 'lo', 'como', 'más', 'pero', 'sus', 'le', 'ya', 'o', 'este', 'sí', 'por

# Stemming{

< Reducir palabras que comparten el mismo significado a su raiz.  
Es simple y rapido>



Las raices no siempre son palabras que si existan

# Stemming{

< Reducir palabras que comparten el mismo significado a una base >

```
import nltk
stemmer = nltk.stem.PorterStemmer()
print(stemmer.stem("improve"))
print(stemmer.stem("improving"))
print(stemmer.stem("improvement"))
print(stemmer.stem("improved"))
print(stemmer.stem("improver"))
```

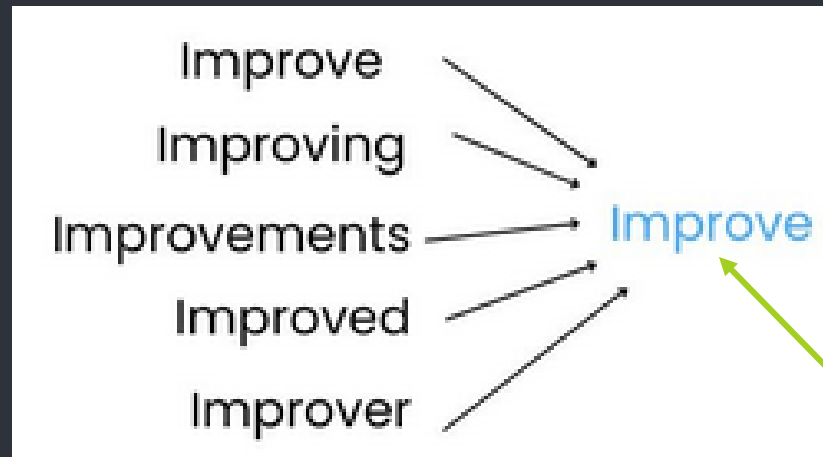
✓ 0.0s

```
improv
improv
improv
improv
improv
```

}

# Lemmetization{

< Reducir palabras que comparten el mismo significado a un lemma, mas lento puesto que implica hacer lookups a vocabulario ya establecido>



Los lemmas son palabras que si existen

# Lemmetization{

< Reducir palabras que comparten el mismo significado a un lemma, mas lento puesto que implica hacer lookups a vocabulario ya establecido>

```
import nltk
nltk.download('wordnet')
from nltk.stem.wordnet import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
print(lemmatizer.lemmatize("improve"))
print(lemmatizer.lemmatize("improving"))
print(lemmatizer.lemmatize("improvement"))
print(lemmatizer.lemmatize("improved"))
print(lemmatizer.lemmatize("improver"))
```

✓ 0.0s

```
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\ivani\AppData\Roaming\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
improve
improving
improvement
improved
improver
```

}



## }



1  
2  
3  
4  
5  
6  
7  
8  
9  
1  
0  
1  
2  
3  
4

Q&A