

INTRODUCCIÓN AL PROCESAMIENTO DEL LENGUAJE NATURAL (PLN)

Leticia Martín-Fuertes Moreno
Lingüista computacional en Bitext
@nimbusaeta
nimbusaeta@gmail.com



Lingwars @ Cylicon Valley

9 de noviembre de 2017

¿QUÉ ES EL PLN?

- Lingüística computacional
- Procesamiento del lenguaje natural (PLN o NLProc)
- Text mining

Área multidisciplinar que combina:

- Lingüística: fonética, sintaxis, semántica
- Informática: programación, aprendizaje automático
- Estadística, probabilidad, análisis de datos
- Lógica, formalización y representación del conocimiento



¿QUÉ ES PLN? PLN ERES TÚ

how old is gary oldman

Todo Imágenes Noticias Videos Shopping Más Configuración Herramientas

Aproximadamente 8.480.000 resultados (0,81 segundos)

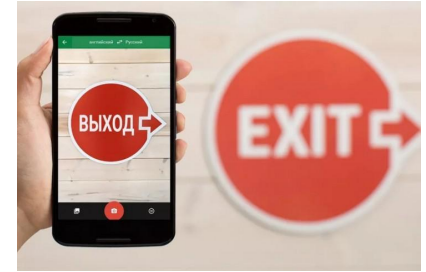
Gary Oldman / Edad

59 años
21 de marzo de 1958



Otras personas también buscan

	Christian Bale 43 años		Alexandra Edenborough 39 años		Michael Caine 84 años
---	---------------------------	---	----------------------------------	---	--------------------------



1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

Usando el Portapapeles:

Mueva o copie el texto al portapapeles haciendo clic en el botón Copiar o Cortar del grupo Portapapeles en la ficha Inicio.

También puede pulsar Ctrl + X para cortar o Ctrl + C para copiar. Cuando hace estas acciones, el texto es movido o copiado en un contenedor electrónico llamado el Portapapeles. Para pegar el texto, simplemente haga clic en el comando Pegar o pulse Ctrl + V.

Nota: Las opciones Copiar, Cortar o Pegar en el menú contextual, son parte de los comandos del Portapapeles.

Ortografía

También

Qmitir Qmitir todas Agregar

También

Cambiar CAMBIAR TODO

También

1. indica que lo expresado a continuación queda incluido en una afirmación precedente
2. señala que la información expresada se



HISTORIA DEL PLN EN UN VISTAZO

- En los 50, al principio de la Guerra Fría, el objetivo del PLN era la **traducción automática**.
- En los 60 aparecen **sistemas expertos** que asistían en la toma de decisiones: sistemas de diálogo que trataban de imitar conversaciones humanas, creación de ontologías para capturar conocimiento del mundo.
- Hasta los 80, la mayor parte de los sistemas de PLN estaban basados en conocimiento y manejaban complejas reglas diseñadas a mano. Influencia de la **lingüística generativa** de Chomsky.
- A partir de esa década, irrumpen las **aproximaciones estadísticas** basadas en sistemas de aprendizaje automático (*machine learning*), que requieren grandes colecciones de datos anotados manualmente. Desarrollo paralelo al aumento de potencia de los ordenadores.
- Actualmente, vivimos un auge de los **sistemas de aprendizaje automático no supervisados** (es decir, no anotados), con especial énfasis en el uso de la Web. Explosión de datos en formato electrónico.
- En la década de 2010 hemos visto el resurgir de los sistemas que utilizan **redes neuronales**.

TAREAS PROPIAS DEL PLN

Tokenización

Dividir un texto en palabras.

```
In [2]: import nltk
        from nltk import word_tokenize, pos_tag
```

```
tokens = word_tokenize(raw)
```

```
print(tokens[:60])
```

```
print(len(tokens))
```

```
[['', 'Emma', 'by', 'Jane', 'Austen', '1816', ''], 'VOLUME', 'I', 'CHAPTER', 'I', 'Emma', 'Woodhouse',  
'', 'handsome', ',', 'clever', ',', 'and', 'rich', ',', 'with', 'a', 'comfortable', 'home', 'and', 'happ  
y', 'disposition', ',', 'seemed', 'to', 'unite', 'some', 'of', 'the', 'best', 'blessings', 'of', 'existenc  
e', ';', 'and', 'had', 'lived', 'nearly', 'twenty-one', 'years', 'in', 'the', 'world', 'with', 'very', 'li  
ttle', 'to', 'distress', 'or', 'vex', 'her', '.', 'She', 'was']
```

```
191673
```

TAREAS PROPIAS DEL PLN

Part-Of-Speech tagging (POS-tagging o etiquetado morfológico)

Asignar a cada palabra su categoría gramatical (o *parte del discurso*). Si puede ser la específica para ese contexto, mejor.

You can't do your work like this. The computer doesn't work!

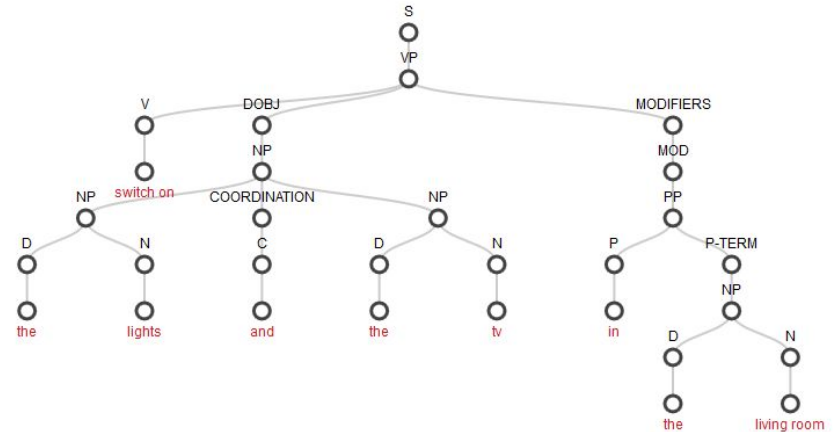
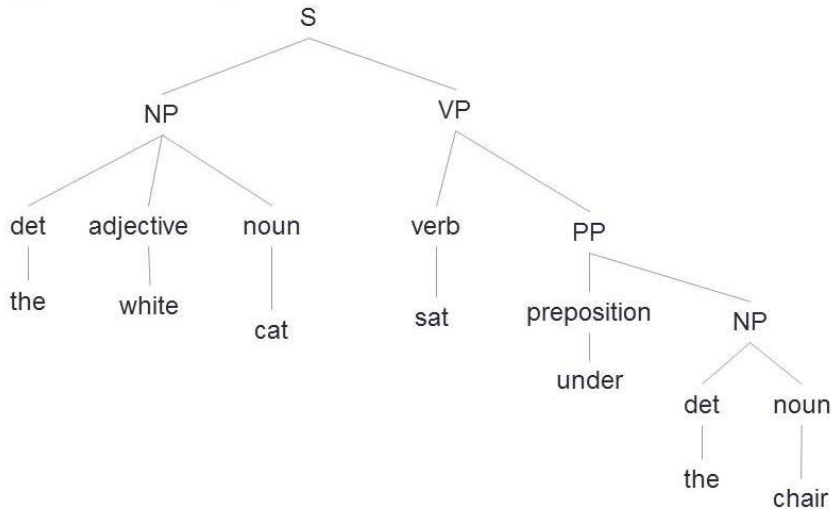
Form	POS
You	pronoun
can	verb
not	adverb
do	verb
your	determiner
work	noun
like	preposition
this	pronoun

Form	POS
The	determiner
computer	noun
does	verb
not	adverb
work	verb

TAREAS PROPIAS DEL PLN

Parsing (o análisis sintáctico)

Analizar una cadena de caracteres según las reglas de una gramática formal. En la práctica, consiste en dividir un texto en constituyentes, obteniendo como resultado un árbol que muestre las relaciones entre estos.



TAREAS PROPIAS DEL PLN

Reconocimiento de entidades (o NER, *Named entities recognition*)

En un texto dado, localizar y clasificar entidades con nombre en categorías definidas, como pueden ser nombres de personas, organizaciones, lugares, expresiones de tiempo, cantidades, valores monetarios, porcentajes, etc.

```
In [15]: Dracula = read_lines("data/corpus_misc/en/literature/Dracula.txt", 200, 500)
sample = Dracula
sample = sample.replace('\n', ' ')
sentences = nltk.sent_tokenize(sample)
tokenized_sentences = [nltk.word_tokenize(sentence) for sentence in sentences]
tagged_sentences = [nltk.pos_tag(sentence) for sentence in tokenized_sentences]
chunked_sentences = nltk.ne_chunk_sents(tagged_sentences, binary=True)

def extract_entity_names(t):
    _entity_names = []
    if hasattr(t, 'label') and t.label():
        if t.label() == 'NE':
            _entity_names.append(' '.join([child[0] for child in t]))
        else:
            for child in t:
                _entity_names.extend(extract_entity_names(child))
    return _entity_names

entity_names = []
for tree in chunked_sentences:
    | entity_names.extend(extract_entity_names(tree))

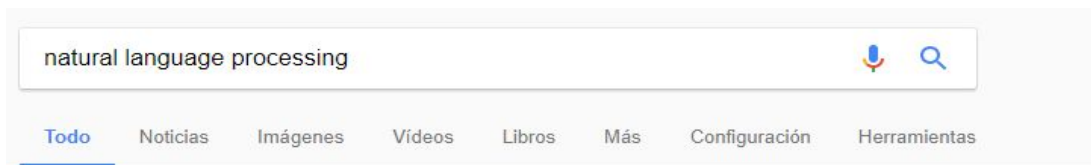
entity_names = set(entity_names)
print(entity_names)

{'Hospadars', 'German', 'Dacians', 'Mittel Land', 'Servian', 'Carpathian', 'Castle Dracula', 'Pass', 'Borgo Prund', 'Germany',
 'English Churchman', 'Wallachs', 'Hotel Royale', 'Attila', 'Borgo Pass', 'Jonathan Harker', 'Golden Mediasch', 'Turkish', 'Kla
 useburgh', 'Turk', 'Count Dracula', 'London', 'St. George', 'Europe', 'Slovak', 'France', 'British Museum', 'English', 'Mina',
 'Golden Krone Hotel', 'Isten', 'Transylvania', 'Moldavia', 'Herr Englishman', 'Danube', 'Cszezs', 'Bukovina', 'Bistritz', 'Sout
 h', 'Vienna', 'East', 'West', 'DRACULA', 'North', 'Huns', 'China', 'Magyars', 'Oriental', '_3 May', 'Slovaks', 'Szekelys', 'Mun
 ich'}
```


TAREAS PROPIAS DEL PLN

Recuperación de información (*information retrieval*)

Devolver una lista de documentos relevantes para una búsqueda.



A screenshot of a search engine interface. The search bar contains the text "natural language processing". To the right of the search bar are icons for voice search and a magnifying glass. Below the search bar is a horizontal menu with tabs: "Todo" (underlined), "Noticias", "Imágenes", "Videos", "Libros", "Más", "Configuración", and "Herramientas".

Aproximadamente 26.700.000 resultados (0,59 segundos)

[Natural language processing - Wikipedia](#)

https://en.wikipedia.org/wiki/Natural_language_processing ▼ Traducir esta página

Natural language processing (NLP) is a field of computer science, artificial intelligence and computational linguistics concerned with the interactions between ...

[Corpus linguistics](#) · [Natural language understanding](#) · [Computational linguistics](#)

[Introduction to Natural Language Processing \(NLP\) 2016 - Algorithmia](#)

<https://blog.algorithmia.com/introduction-natural-language-proce...> ▼ Traducir esta página

11 ago. 2016 - **Natural Language Processing**, or **NLP** for short, is a field of study focused on the interactions between human language and computers.

[The Definitive Guide to Natural Language Processing - MonkeyLearn](#)

<https://monkeylearn.com/.../definitive-guide-natural-language-pr...> ▼ Traducir esta página


29 oct. 2015 - A guide that gives an introduction to **Natural Language Processing**, explaining how can a machine understand text, important concepts and ...

TAREAS PROPIAS DEL PLN

Análisis de opinión (*sentiment analysis*)

Determinar de una forma automática la actitud del hablante hacia un tema, o la polaridad contextual o reacción emocional a un determinado producto, documento, interacción o evento.

Besides, I have had a good customer service experience. The salesman John Faraday was very nice!



Besides, I have had a ¹ good ¹ customer ¹ service ¹ experience. The ² salesman ² John Faraday was ² very nice!

LEGEND color key

SENTIMENT

Sentiment topic

Positive sentiment text

Negative sentiment text

¹ Text and topic link



Language Technology

making good progress

mostly solved

Spam detection

Let's go to Agra!

Buy VIAGRA ...

Part-of-speech (POS) tagging

ADJ ADJ NOUN VERB ADV

Colorless green ideas sleep furiously.

Named entity recognition (NER)

PERSON ORG LOC

Einstein met with UN officials in Princeton

Sentiment analysis

Best roast chicken in San Francisco!

The waiter ignored us for 20 minutes.

Coreference resolution

Carter told Mubarak he shouldn't run again.

Word sense disambiguation (WSD)

I need new batteries for my *mouse*.

Parsing

I can see Alcatraz from the window!

Machine translation (MT)

第13届上海国际电影节开幕...

The 13th Shanghai International Film Festival...

Information extraction (IE)

You're invited to our dinner party, Friday May 27 at 8:30



still really hard

Question answering (QA)

Q. How effective is ibuprofen in reducing fever in patients with acute febrile illness?

Paraphrase

XYZ acquired ABC yesterday

ABC has been taken over by XYZ

Summarization

The Dow Jones is up

The S&P500 jumped

Housing prices rose



Economy is good

Dialog

Where is Citizen Kane playing in SF?



Castro Theatre at 7:30. Do you want a ticket?



LO ÚLTIMO DE LO ÚLTIMO

Aplicar *machine learning*, *deep learning* y redes neuronales a PLN para desarrollar bots conversacionales (chatbots)

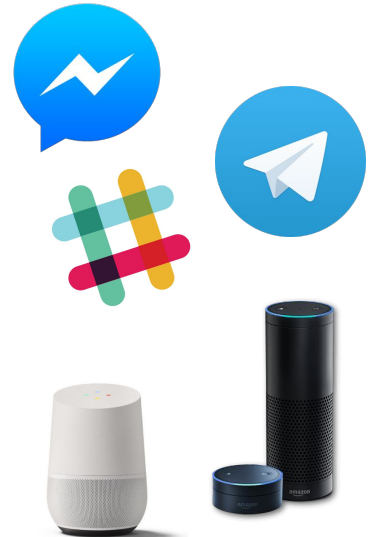
Plataformas PLN



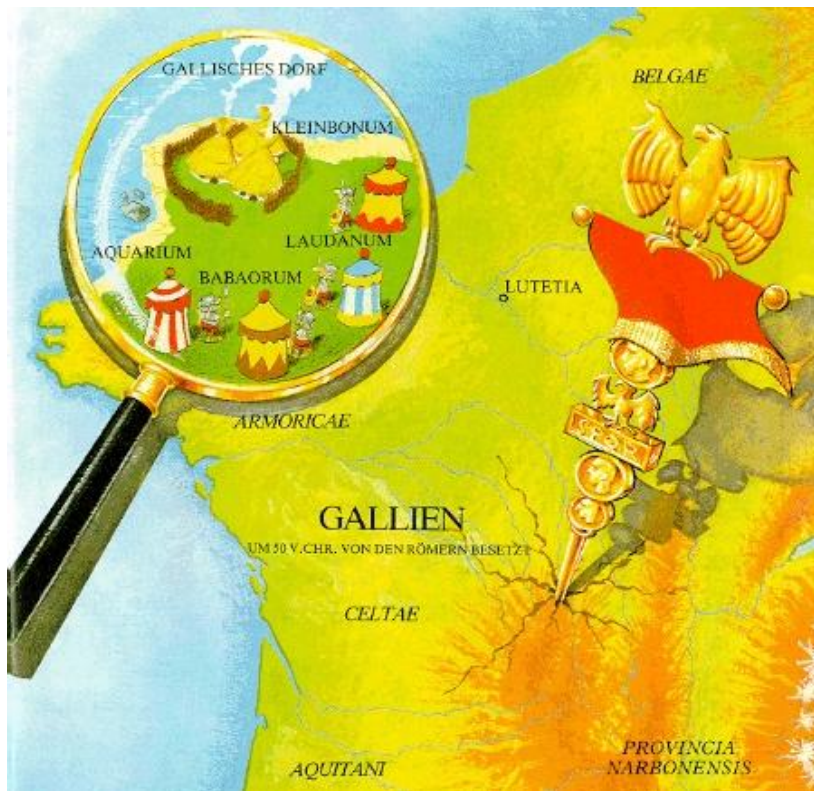
Construcción



Canales



Y EN UNA IRREDUCTIBLE ALDEA GALA...



Lingwarriors
de Lingwars

Charlas y
talleres

GAPLEN

OBJETIVOS

- Lo que nos une es el interés por el PLN. No hay un lenguaje de programación preferido.
- Quién: gente «de letras» + gente «de ciencias»
- Se valoran todos los conocimientos y todos los niveles
- Enfoque lingüístico sobre el enfoque estadístico

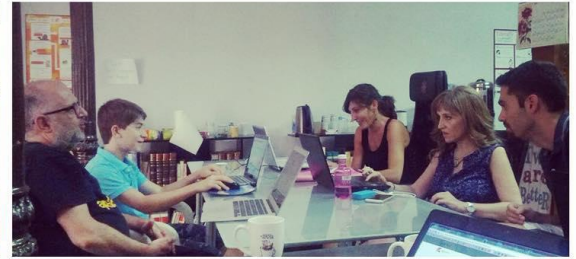


TRAYECTORIA

22/07/2015 - Primera quedada. Hasta mediados de noviembre: quedadas semanales los miércoles o jueves

03/10/2015 - Jornada entera: **presentación de proyectos** desarrollados en Lingwars + Introducción a Python

14/11/2015 - Jornada entera, edición especial corpus: AntConc + **presentación de proyectos** (Neutrón, Aracne, enclitizador) + Recuperación de información



TRAYECTORIA

09-23/04/2016 - Taller de **tratamiento de datos digitales** (scraping, análisis de corpus, visualización de datos) por Javi G. Sogo + Tania Karaseva



07/05/2016 - Taller de la **API de Twitter** por Javi G. Sogo

09/07/2016 - Charla sobre **semántica estructural** por Mónica González Manzano

16/07/2016 - Taller de **Jupyter para lingüistas** por Javi G. Sogo

12/12/2016 - Charla sobre la **sílaba y Prolog** por Elena Álvarez Mellado



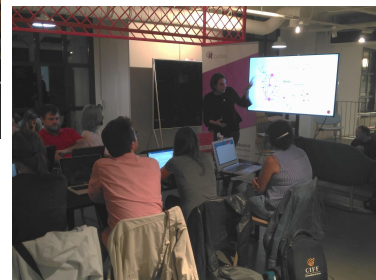
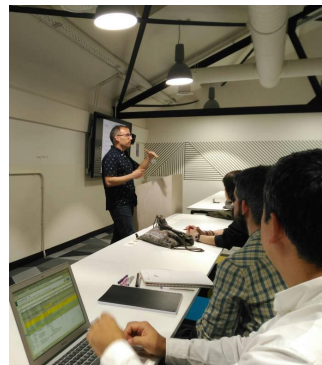
TRAYECTORIA

22/02/2017 - Charla sobre **cómo hacerse una carrera en PLN** por Víctor Peinado

08/07/2017 - Taller de **NLTK** por Leticia Martín-Fuertes, Tania Karaseva y Javi G. Sogo

05/10/2017 - GAPLEN #1 - Charla de **introducción a PLN** por Luis Anke Espinosa

22/10/2017 - Charlas sobre **ontologías + PLN en R** por Lorena Giusio y Claudia Guirao @ R-Ladies Madrid



TRAYECTORIA

02/11/2017 - GAPLEN #2 - Charla **Rapidminer** por Enrique Puertas

09/11/2017 - Charlas **Corpus + Ontologías** por Leticia Martín-Fuertes y Javi G. Sogo @ UVa

09/11/2017 - Charlas de **introducción a PLN + ontologías** por Leticia Martín-Fuertes y Javi G. Sogo @ Cylicon Valley

22/11/2017 - Charla sobre **resumen automático** por Luis Anke Espinosa



???

ESPÍRITU

sobre los embeddings (skip-gram etc) la verdad podíamos ver d hablar en algún momento... es un tema ad+ d ser muy hot hoy, tiene mucha chicha lingüística 9:48 AM

A mi si me gustaría aprender algo más de los embeddings, aunque tengo mucho que aprender todavía 9:51 AM

Yo probablemente también me apunte :) 9:56 AM

guay! 9:57 AM

Tarde de trueque. Yo os echo una mano con la configuración del ordenador y a mi me habláis de temas lingüísticos; embeddings, skip-gram, ... 9:57 AM

xD bueno se hará lo q se pueda 9:57 AM

¡o qué guay, esa ha sido siempre la idea :) 9:57 AM

LingWars @lingwars · 8 May 2016
El currazo que se ha pegado @jsgogo últimamente nos ha dejado 2 magníficos #tutoriales: cómo scrapear la Web lingwars.github.io/blog/scrape-xp...
Translate from Spanish
1 2 4

LingWars @lingwars Following

Y cómo conectarse a la API de Twitter lingwars.github.io/blog/twitter-s... (sí, contiene información sobre qué leches es una API) #tutoriales



juliverx @GoogleCampusMa... Following

juliverx Aquí estamos: lingüinis y techies con @nimbusaeta y @lirondos en @lingwars playmobiles Mola. Me apunto a la próxima. Sea lo que sea 😊
elizabuben Love your pictures!
nimbusaeta @juliverx ¡Gracias!
@playmobiles Somos una especie de grupo de apoyo para aprender procesamiento de lenguaje natural, resumidamente :)

LingWars @lingwars

Following

«Yo estoy emocionado solo con el hecho de que estemos aquí reunidas unas cuantas personas hablando de estas cosas» se ha dicho hoy aquí <3

Elena A. Mellado @lirondos · 19 ago. 2015
En respuesta a @lirondos
Representación gráfica de lingufriendly.



Elena A. Mellado @lirondos · 6 Dec 2016
"Dime, oh Azrael, en qué idioma está escrito este texto".
Hoy ha sido un día provechoso :)_
Translate from Spanish
:ta aszetika deritzon a: (1643) euskal prosa
o Agerre Azpilkueta (A. 3333333333333333
araldeko Saran ema 2222222222222222
, geroago lapurtera klz. 66666666666667
azpideak bilatzen ditu. 7.77777777777778
enetako sortzailea, Le

Elena A. Mellado @lirondos Following

De esta criatura y del cómo se hizo hablaremos entre otras cosas el lunes. Es gratis, es chachi, somos majos

Leticia MFM @nimbusaeta

Bieeeeeen 🙌🙌 ya hemos llegado al momento típico en que alguien dice que las regex son un quebradero de cabeza xD @claudiaguirao

¿QUÉ NOS DEPARA EL FUTURO?



THAT'S ALL, FOLKS!



lingwars@gmail.com

[@lingwars](#)