



DataSpartan



LLMs: La próxima revolución de la inteligencia artificial

¿Qué son los LLMs? ¿Cómo funcionan? ¿Cómo puedo sacarles más partido?



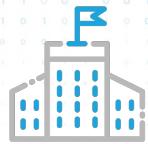
[https://github.com / rorisDS / workshop_llms](https://github.com/rorisDS/workshop_llms)





DataSpartan

Expertos en las áreas de Inteligencia Artificial (IA), Big Data, Investigación Cuantitativa, Computación de Alto Rendimiento (HPC), Machine Learning (ML) y Modelización de Riesgos.



PYMES

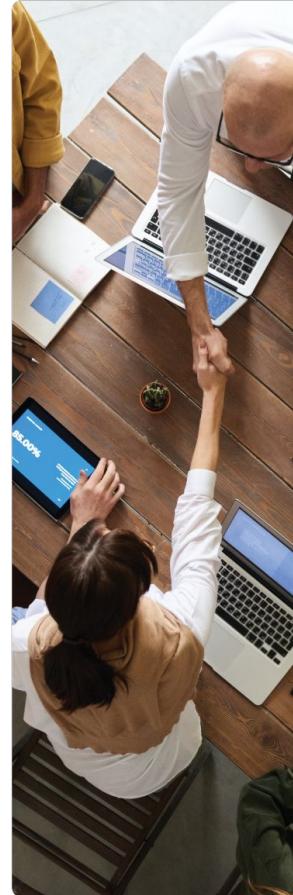


EMPRENDEDORES



EQUIPOS I+D

Investigamos, diseñamos, construimos, implementamos y validamos soluciones, para facilitar que las empresas puedan alcanzar su potencial, de forma más rápida y asequible de lo que sería posible contando únicamente con sus propios medios.



Víctor Manuel Alonso Rorís

- Ingeniería de Telecomunicaciones (2003 - 2009)
- Doctor en Telemática (2017)
 - *Tecnologías Semánticas - Representación de conocimiento*
 - *34 publicaciones científicas (2011-2019)*
- Investigador en la Universidad de Vigo como investigador (9 años)
- Data Scientist en DataSpartan (6 años)
 - *Machine Learning Engineer – NLP specialist*



3 Workshop en el Foro de Empleo en la UVigo



Contenidos

1. Conceptos básicos de Machine Learning
2. ¿Qué es NLP?
3. Introducción a los LLMs
4. Cómo sacar partido a los LLMs

Advertencia: Muchas explicaciones pueden no ajustarse a una definición exacta o real pero se ha hecho así para facilitar la comprensión del lector.

Alternativa 1 de Gemini: Es importante destacar que algunas de las explicaciones proporcionadas podrían no ajustarse a una definición precisa o completamente real. Se han simplificado con el objetivo de facilitar la comprensión del lector.

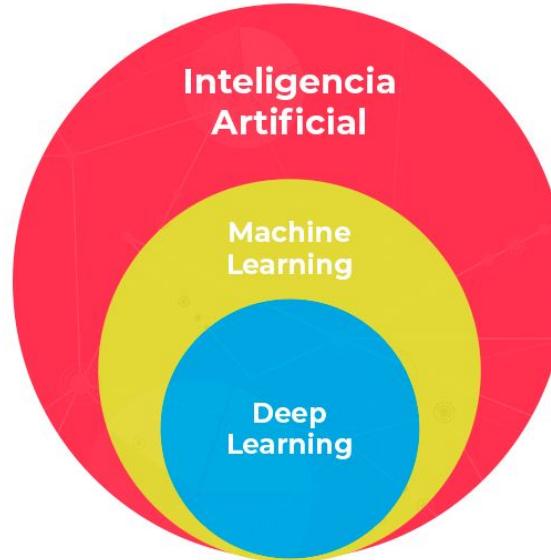
Alternativa 2 de Gemini: Algunas explicaciones podrían no ser completamente precisas o ajustarse a la realidad. Se han simplificado para facilitar la comprensión del lector.



Conceptos básicos de Machine Learning



Definiciones



IA: Combinación de algoritmos planteados con el propósito de crear máquinas que presenten las mismas capacidades que el ser humano.

Machine Learning: Rama de la Inteligencia artificial (IA) que estudia como dotar a las máquinas de capacidad de aprendizaje

Deep Learning: algoritmo automático jerárquico que emula el aprendizaje humano con el fin de obtener ciertos conocimientos.

[source](#)

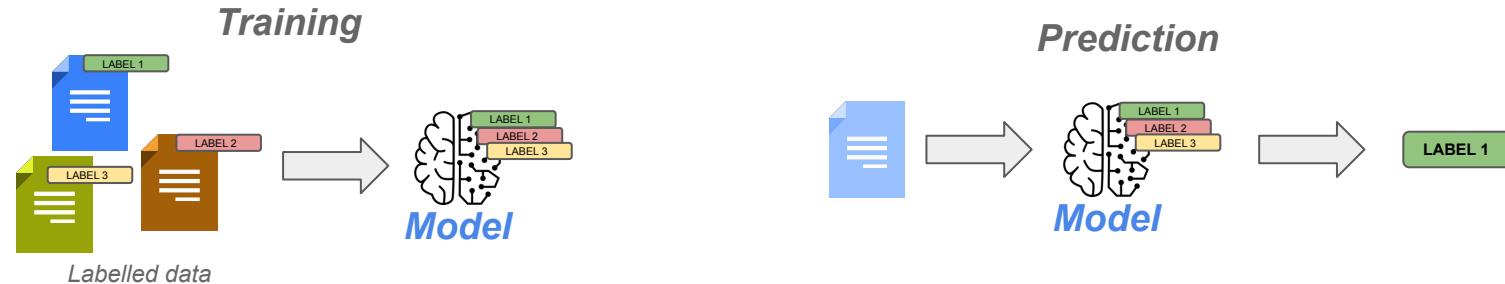


Supervised Learning

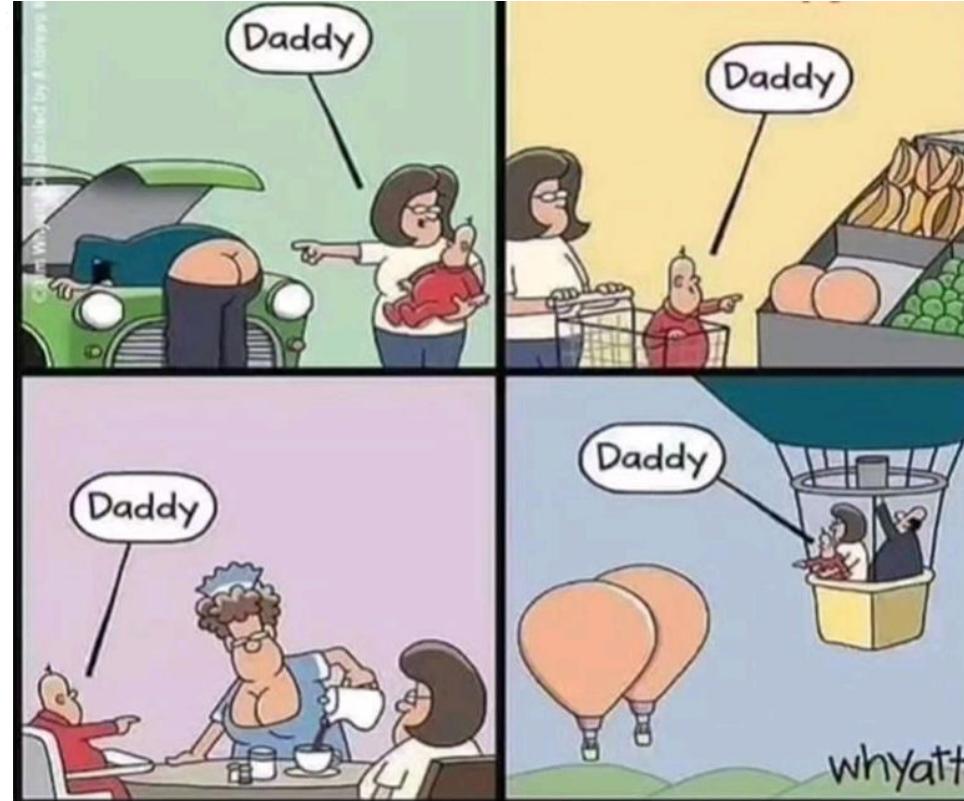
Tipo de ML que utiliza datos etiquetados para entrenar models ML.

Datos etiquetados significa que la salida esperada es conocida.

Los algoritmos miden su precisión a través del error de predicción, ajustándose hasta que el error sea suficientemente minimizado.



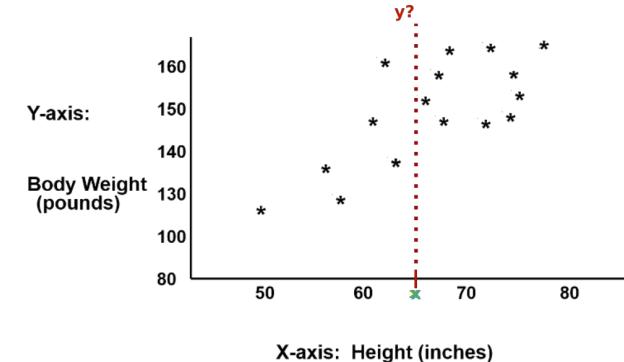
Supervised Learning



Simple Linear Regression

Predecir un simple output (y) a traves de 1-d input (x)

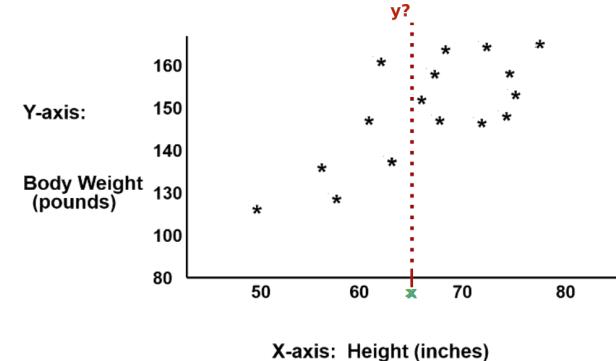
Por ejemplo, predecir el peso corporal (y) por la altura (x)



Simple Linear Regression

Predecir un simple output (y) a través de 1-d input (x)

Por ejemplo, predecir el peso corporal (y) por la altura (x)

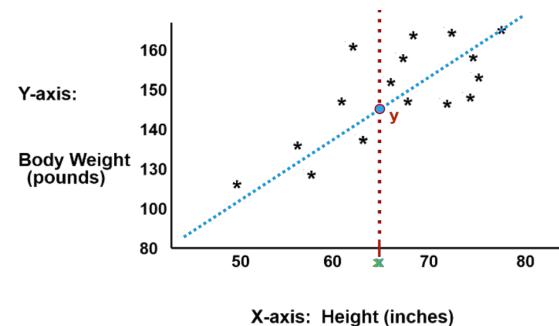


Formula de la recta:

$$h(x) = m \cdot x + b = y'$$

Donde:

- m es la pendiente
- b es el término independiente



Simple Linear Regression

1. Valor "aleatorio" para m y b
2. Introduce x en la función \Rightarrow predice y'
3. Calcula error entre y' y el y esperado
4. Ajusta m y b en base al error
 \Rightarrow Descenso del gradiente
5. Repite el proceso múltiples veces para cada ejemplo del dataset de entrenamiento

$$h(x) = m*x + b = y'$$

$$h_0(x) = 2.3*x + 1 = y'$$

$$h_1(x) = 4.1*x + 5 = y'$$

...

$$h_M(x) = 3.032*x + 2.95 = y'$$



Simple Linear Regression

1. Valor "aleatorio" para m y b
2. Introduce x en la función \Rightarrow predice y'
3. Calcula error entre y' y el y esperado
4. Ajusta m y b en base al error
 \Rightarrow Descenso del gradiente
5. Repite el proceso múltiples veces para cada ejemplo del dataset de entrenamiento

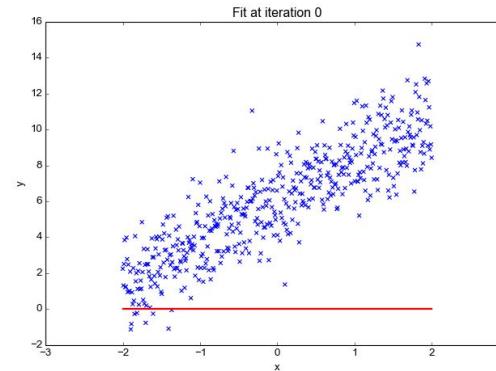
$$h(\mathbf{x}) = m^*\mathbf{x} + b = y'$$

$$h_0(x) = 2.3*x + 1 = y'$$

$$h_1(x) = 4.1*x + 5 = y'$$

...

$$h_M(x) = 3.032*x + 2.95 = y'$$



Multiple Linear Regression

Vector n-dimensional de input : $\mathbf{x} = (x^1, x^2, x^3, \dots, x^N)$

$$h(\mathbf{x}) = \theta_0 * x_0 + \theta_1 * x_1 + \dots + \theta_N * x_N + b = y'$$



Multiple Linear Regression

Vector n-dimensional de input : $\mathbf{x} = (x^1, x^2, x^3, \dots, x^N)$

Output m-dimensional : $\mathbf{y} = (y^1, y^2, y^3, \dots, y^M)$

$$h_1(\mathbf{x}) = \theta_0^1 * x_0 + \theta_1^1 * x_1 + \dots + \theta_N^1 * x_N + b^1 = y_1,$$

$$h_2(\mathbf{x}) = \theta_0^2 * x_0 + \theta_1^2 * x_1 + \dots + \theta_N^2 * x_N + b^2 = y_2,$$

$$h_3(\mathbf{x}) = \theta_0^3 * x_0 + \theta_1^3 * x_1 + \dots + \theta_N^3 * x_N + b^3 = y_3,$$

...

$$h_M(\mathbf{x}) = \theta_0^M * x_0 + \theta_1^M * x_1 + \dots + \theta_N^M * x_N + b^M = y_M,$$



Multiple Linear Regression

Vector n-dimensional de input : $\mathbf{x} = (x^1, x^2, x^3, \dots, x^N)$

Output m-dimensional : $\mathbf{y} = (y^1, y^2, y^3, \dots, y^M)$

$$h_1(\mathbf{x}) = \theta^1_0 * x_0 + \theta^1_1 * x_1 + \dots + \theta^1_N * x_N + b^1 = y_1,$$

$$h_2(\mathbf{x}) = \theta^2_0 * x_0 + \theta^2_1 * x_1 + \dots + \theta^2_N * x_N + b^2 = y_2,$$

$$h_3(\mathbf{x}) = \theta^3_0 * x_0 + \theta^3_1 * x_1 + \dots + \theta^3_N * x_N + b^3 = y_3,$$

...

$$h_M(\mathbf{x}) = \theta^M_0 * x_0 + \theta^M_1 * x_1 + \dots + \theta^M_N * x_N + b^M = y_M,$$

$$\left[\begin{array}{cccc} x_0 & x_1 & x_2 & \dots & x_N \end{array} \right]$$

$$\begin{matrix} \theta^1_0 & \theta^1_1 & \dots & \theta^1_N \\ \theta^2_0 & \theta^2_1 & \dots & \theta^2_N \\ \theta^3_0 & \theta^3_1 & \dots & \theta^3_N \\ \dots \\ \theta^M_0 & \theta^M_1 & \dots & \theta^M_N \end{matrix}$$

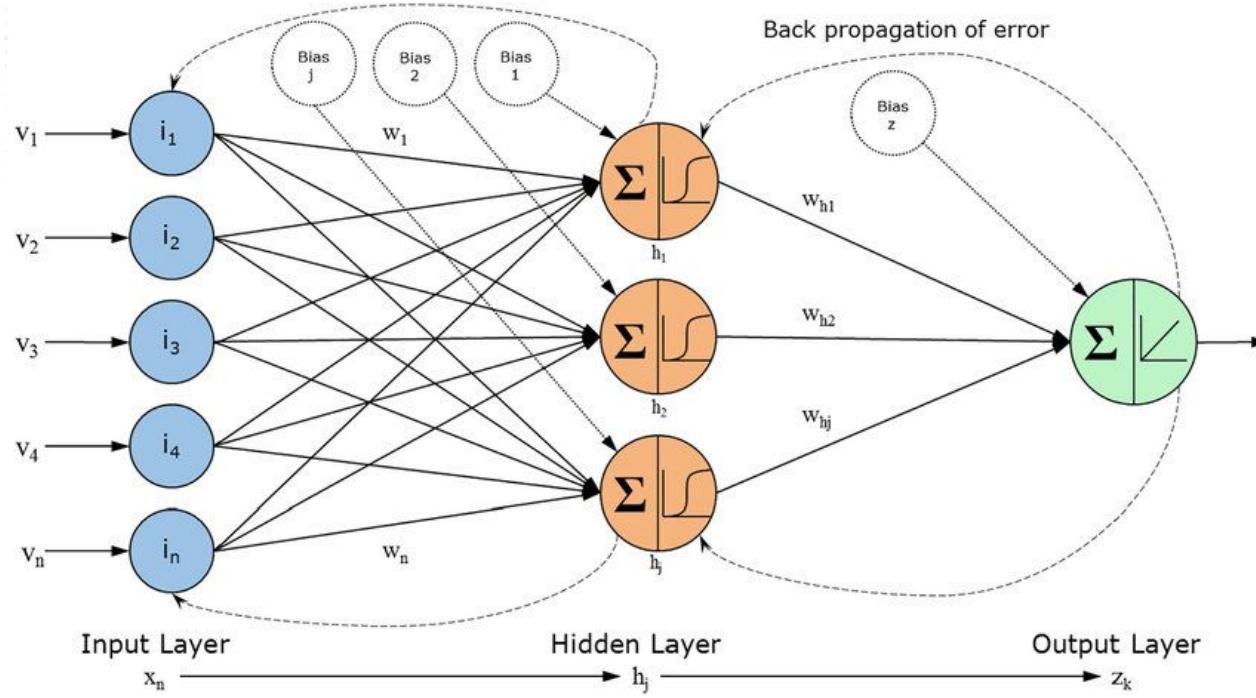


$$\begin{matrix} b_0 \\ b_1 \\ b_2 \\ \dots \\ b_N \end{matrix}$$

OPERACIONES CON MATRICES



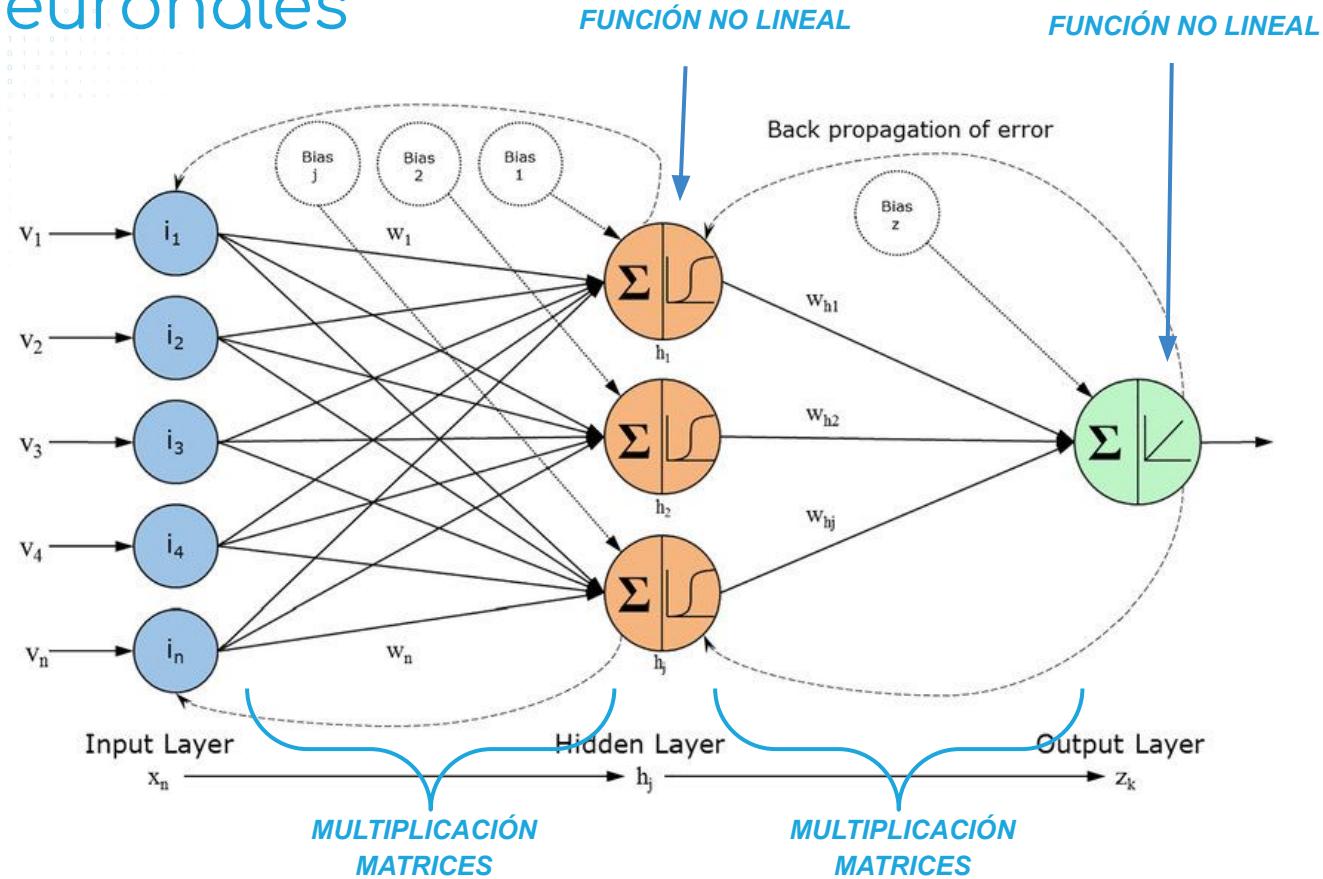
Redes Neuronales



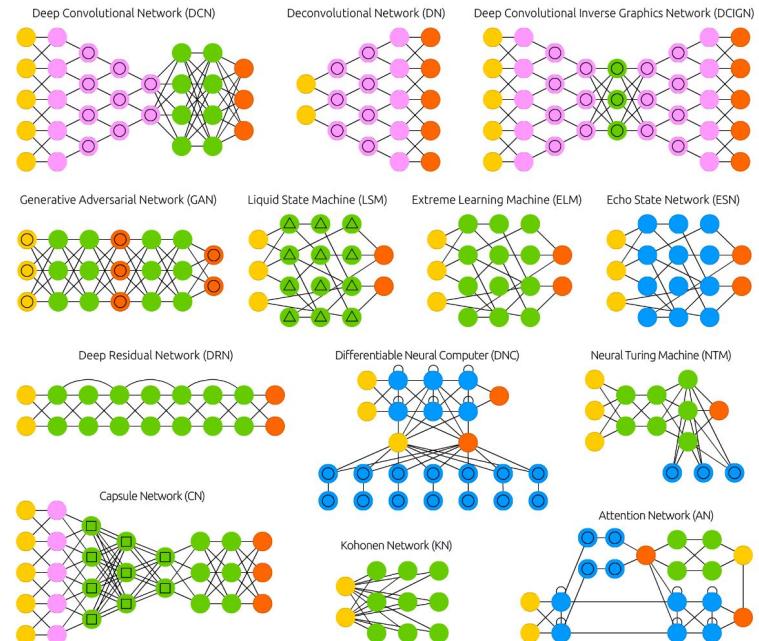
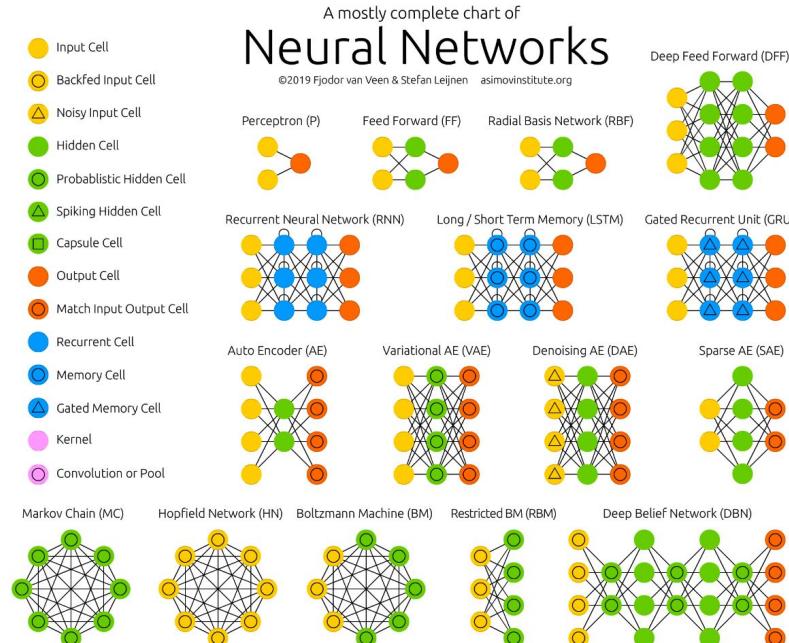
[source](#)



Redes Neuronales



Arquitecturas de Redes Neuronales



[source](#)



Conceptos Machine Learning

Operaciones matemáticas ⇒ Identificación de patrones

No hay conciencia! Lo siento Sarah Connor!

Comportamiento modelado por datos de entrenamiento

Modelos racistas o sexistas ⇒ Tus datos apestan!

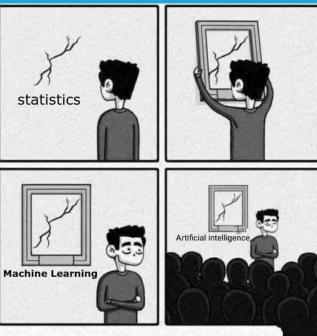
Librerías abstractas matemáticas: [Pytorch](#), [Tensorflow](#), [Keras](#), [Theano](#), ...

Gente mucho más inteligente que yo!

Quieres aprender más:

- [Curso de Coursera](#)
- [Curso de Google](#)
- [Curso de FastAI](#)

Me patiently waiting for AI to begin its war against humanity



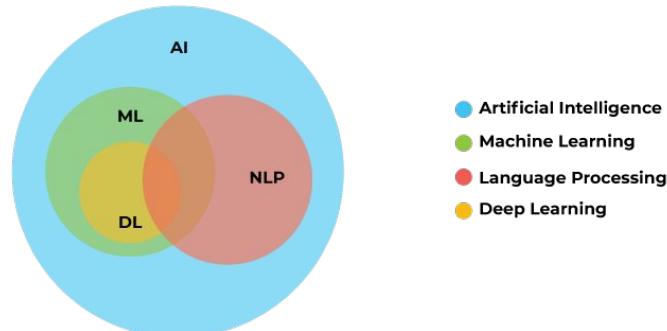
¿Qué es NLP?



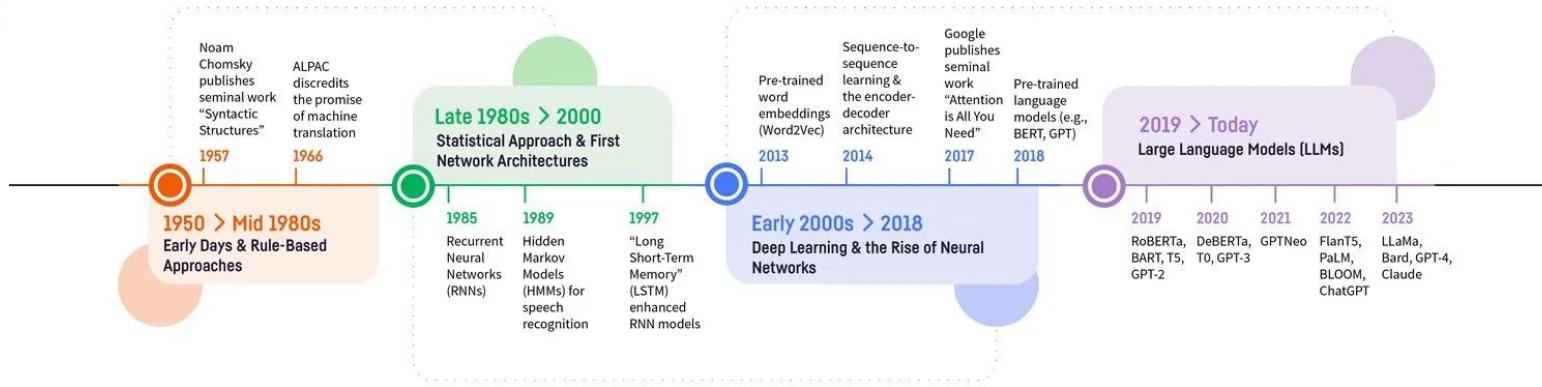
Natural Language Processing

Lenguaje Natural es la lengua o idioma hablado o escrito por seres humanos para propósitos generales de comunicación ⇒ Ambiguo y desestructurado

Procesado del Lenguaje Natural (PLN), **Natural Language Processing (NLP)** en inglés, es la rama de la IA encargada de dar a los ordenadores la capacidad de comprender textos en lenguaje natural.

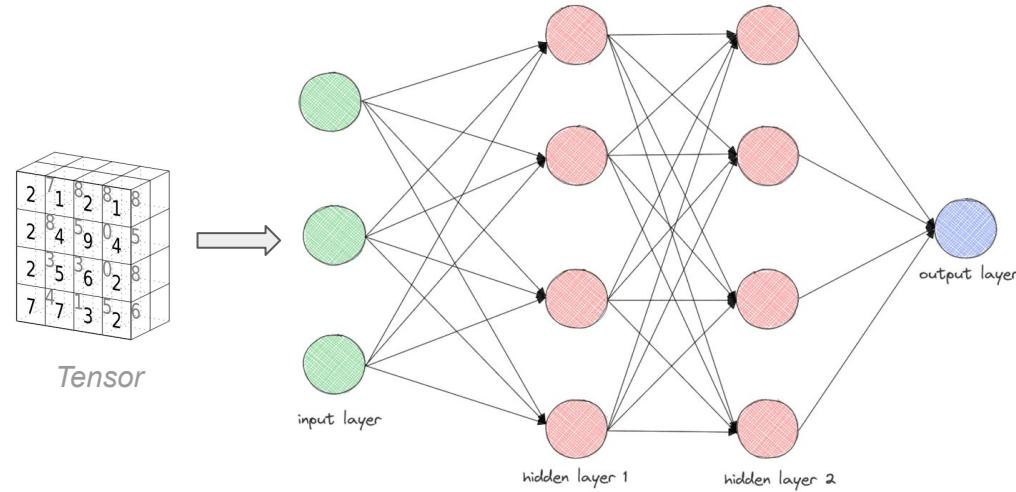


The history of NLP



Neural NLP

Los modelos ML usan tensores como input ⇒ Vectores numéricos

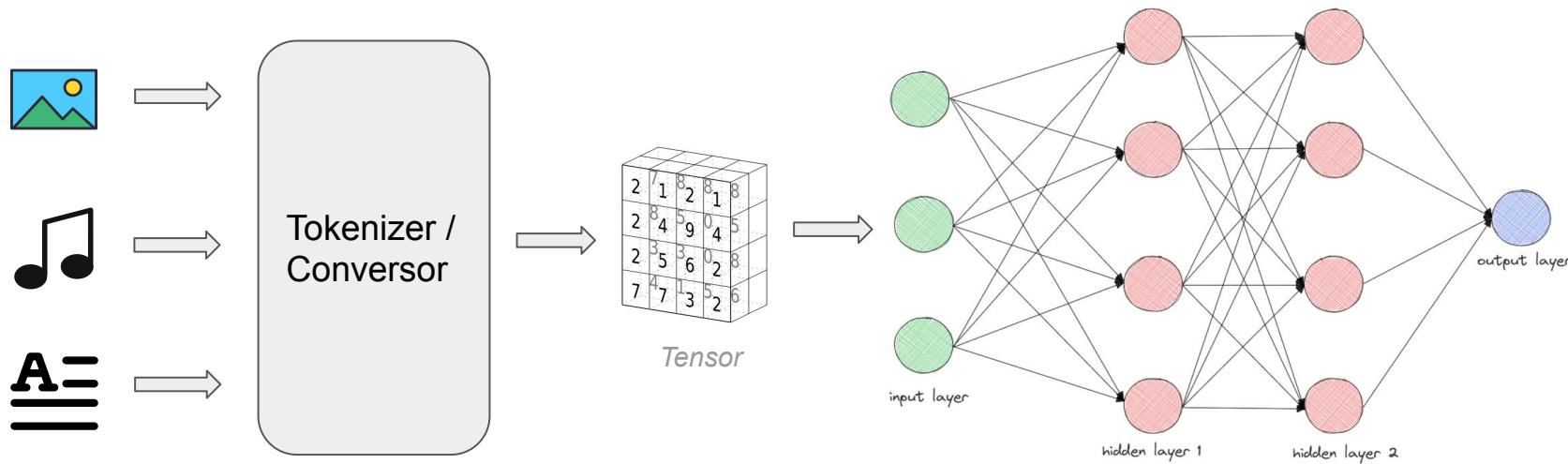


Neural NLP

Los modelos ML usan tensores como input ⇒ Vectores numéricos

¿Qué pasa con los modelos que toman como input: Imágenes? Audio? **Texto?**

⇒ Antes es necesario convertir el input a tensor (features)



Embeddings

Representación numérica de palabras o frases ⇒ vectores

Word	Embedding			
ábaco	0.19	-0.32	...	0.75
abad	0.41	-1.27	...	-0.06
...			...	
zueco	-0.76	-1.09		1.35

¿Tienen estos embeddings algún sentido?



Matemáticas del lenguaje

Semántica de cualquier idioma

⇒ Siguen reglas estadísticas

Técnica Masking

⇒ Introducir ruido en frases y entrenar el modelo para corregirlo

⇒ Dataset Auto-supervisado

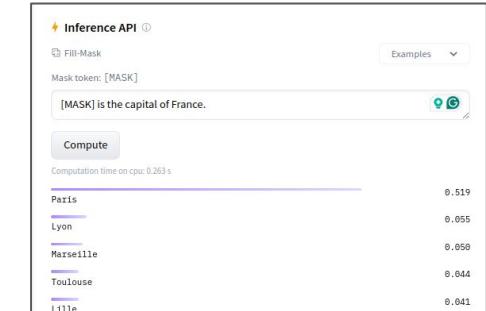
La capital de Francia es [MASK]

[MASK], capital de Francia, es conocida como la ciudad de las luces

Madrid es a España lo que [MASK] es a Francia

En mi visita a Francia, he estado en [MASK] visitando la Torre Eiffel

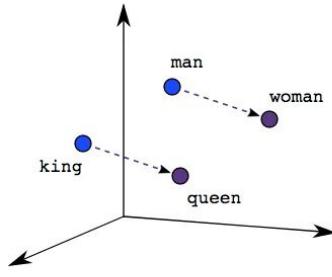
Las grandes capitales europeas, Madrid, Londres, [MASK], Roma y Berlín, acogen ...



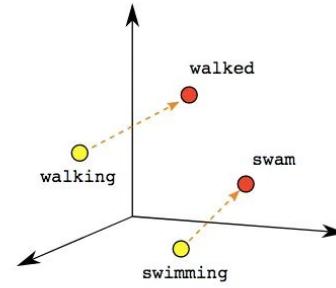
¿Y si fuera un idioma extraterrestre pero dispusieras de millones de textos de ejemplo?



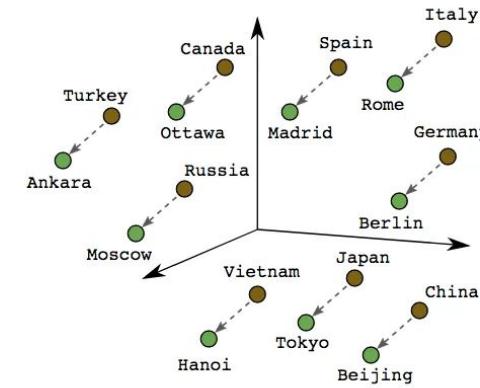
Matemáticas del lenguaje



Male-Female



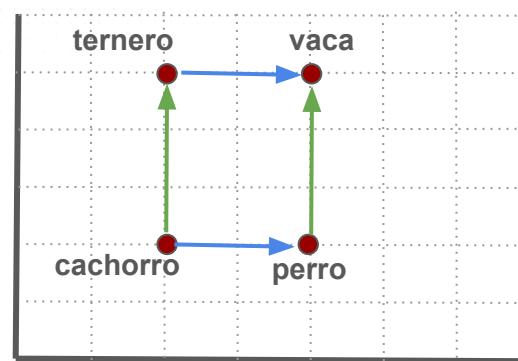
Verb Tense



Country-Capital



Matemáticas del lenguaje



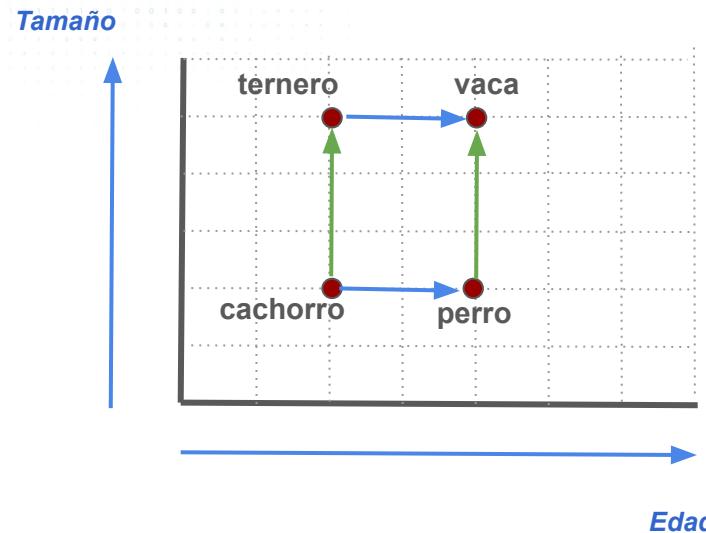
Analogía

Un cachorro es a un perro lo que un ternero a una vaca

Un cachorro es a un ternero lo que un perro a una vaca



Matemáticas del lenguaje



Cada eje representa una propiedad

⇒ Embeddings (vector de multiples dimensiones,

Modelos ML \Rightarrow reajuste automático basado en textos de ejemplos

⇒ Propiedad de cada componente es desconocida

¿Dónde ubicarías una ballena? ¿Y un novillo?



Explosión de NLP

2017 - Google publica "Attention is all you need"

⇒ Se propone la arquitectura ML *Transformers*

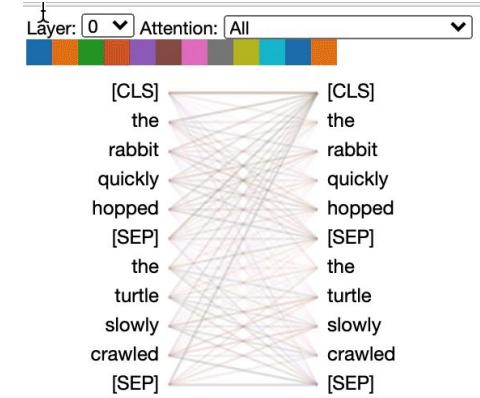
2018 - Pre-trained Language models ⇒ *Fine-tuning*

⇒ *ULMFit, Bert*

⇒ *Datasets de entrenamiento mucho más pequeños*

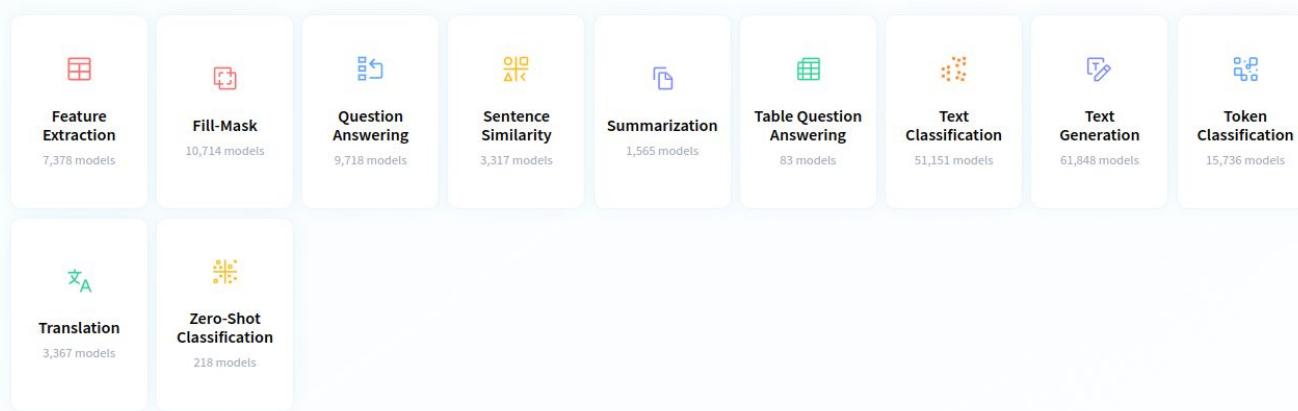
2022 - OpenAi anuncia chatGPT

⇒ *Populariza el término LLM*



Tareas de NLP

Natural Language Processing



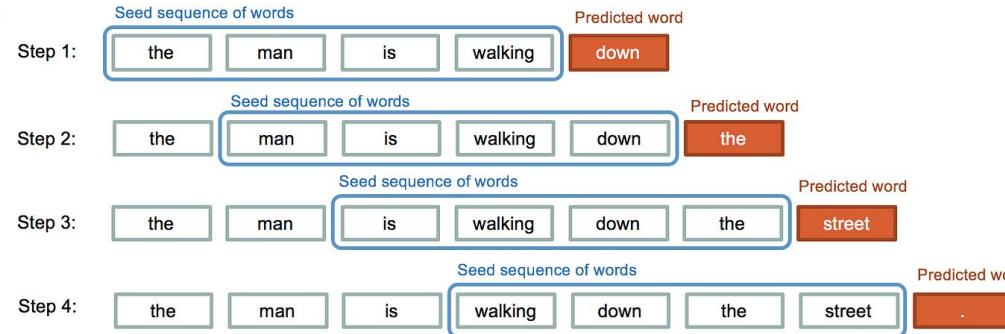
<https://huggingface.co/tasks>



Tareas de NLP: Text Generation

Genera nuevo texto a partir de un texto dado

⇒ predecir la siguiente palabra más probable (o final de texto)



[source](#)

Los datos de entrenamiento son claves para el desempeño del modelo
⇒ No significa que copie textos!

¿Deberían estar los escritores/guionistas/periodistas preocupados? ¿de dónde se sacan los textos de entrenamiento?
¿es lícito usar textos de internet? ¿es lícito reemplazar trabajos por modelos de generación de texto?



Conceptos NLP

Explosión NLP por uso de ML

→ Explotar las relaciones semánticas del lenguaje

Generación de texto ⇒ proceso estadístico

No hay conciencia! Pero puede altar altas cuotas de calidad!

Librerías NLP: [HuggingFace](#), [Spacy](#), [NLTK](#), [Sentence-Transformers](#), ...

Gente mucho más inteligente que yo!

Quieres aprender más:

- [Curso de HuggingFace](#)
- [Curso de Stanford](#)
- [Spacy Tutorial](#)



Tell Me This 20 hours ago (edited)

Human: What do we want?

Computer: Natural language processing!

Human: When do we want it!?

Computer: When do we want what?

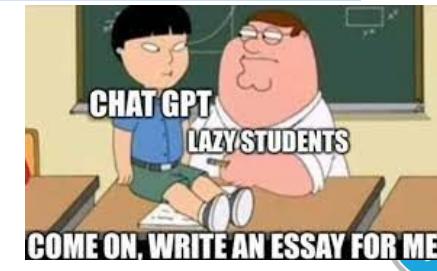
Reply • 203

[View reply](#) ▾



Finetune
BERT

4



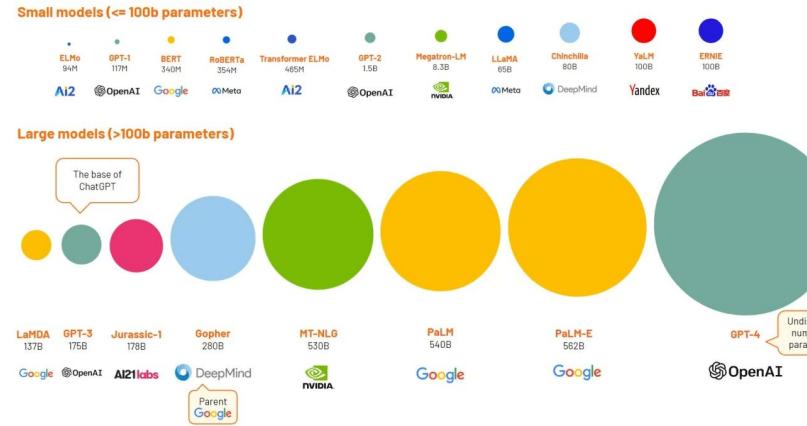
Introducción a los LLMs



Definición LLMs

Large Language Model (**LLM**) es:

- Modelo ML para **NLP** basado en la arquitectura **Transformers**,
- Contiene **millones de parámetros** (conocidos como: m , θ , b , pesos, ...),
- Ha sido entrenado en una **cantidad ingente de textos** (corpus)
- **Propósito general** aunque pueda refinarse para tareas concretas (fine-tuning) de procesamiento de lenguaje natural.



Soluciones LLMs

En la actualidad existen multitud de LLMs disponibles

Soluciones LLM:

- Open-source
 - Mantener control del entorno ⇒ **cuestiones de privacidad**
código, datos de entrenamiento, interacciones con el usuario
 - Requisitos computacionales muy elevados
- Comercial
 - Mejores modelos
 - Facilidad de uso

⇒ acceso vía API (requisitos computacionales trasladados al proveedor)

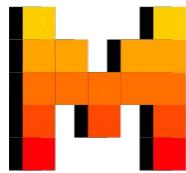




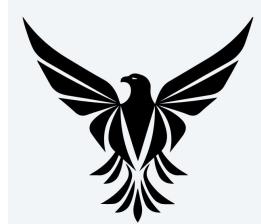
Open-Source LLM



[Llama 2](#)



[MistralAI](#)



[falcon](#)



[BLOOM](#)

**MPT-7B
LLM**

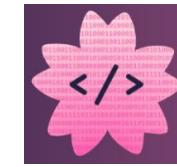
[MPT-7B](#)



[Vicuna](#)



[Tulu](#)



[SantaCoder](#)

Links of interest for fine-tuning an Open-Source LLM:

- [Transformer Reinforcement Learning](#)
- [LoRA](#)
- [Quantization](#)



LLM comerciales



[OpenAI](#)



Conceptos LLM

LLM son modelos enormes
⇒ Demanda recursos computacionales
Tu portatil no esta preparado!

More parameters ⇒ More fun!
⇒ Capacidades sorprendentes
¿Dónde está el límite?

¿De verdad necesitas un LLM?

Librerías LLM: [HuggingFace](#), [LangChain](#)

Quieres aprender más:

- [LLM explained](#)
- [Build a chatGPT from scratch](#)



Cómo sacar partido a los LLMs



Prompt

Los LLM generan texto a partir de un texto de entrada.

Este texto indica el comportamiento esperado por el LLM

⇒ Ese texto es un **PROMPT**

Buen prompt ⇒ Mejores resultados

Prompt Engineering: desarrollar prompts que guíen al LLM para que lleve a cabo tareas especializadas sin ser explicitamente entrenado para ello.

A Actúa como si fueras un experto en astrofísica y explícame qué son los agujeros negros

iPor supuesto! Como experto en astrofísica, puedo explicarte qué son los agujeros negros. Los agujeros negros son regiones extremadamente densas en el espacio-tiempo que se forman cuando una estrella masiva colapsa bajo la influencia de su propia gravedad. La gravedad es la fuerza que atrae a los objetos hacia otros objetos con masa. Cuando una estrella masiva se queda sin combustible nuclear en su núcleo, no puede contrarrestar la fuerza de gravedad que actúa hacia adentro y comienza a colapsar bajo su propio peso. Durante este colapso, la materia de la estrella se comprime en un volumen muy pequeño. La gravedad en el colapso es tan intensa que la estrella colapsada forma un objeto con una fuerza gravitacional tan poderosa que nada, ni siquiera la luz, puede escapar de su atracción. Esto es lo que llamamos un agujero negro. La región central de un agujero negro se llama singularidad, donde toda la masa colapsada se concentra en un punto infinitesimal y la gravedad se vuelve infinita. Alrededor de la singularidad, hay una región conocida como el horizonte de sucesos. Es una frontera imaginaria más allá de la cual la gravedad del agujero negro es tan fuerte que nada puede escapar de su atracción, incluida la luz. El horizonte de sucesos



Reglas de oro para el uso de prompts

Sé concreto, descriptivo y lo mas detallado posible sobre la salida esperada (tamaño, formato, estilo, etc.)

Less effective ✗:

Write a poem about OpenAI.

Better ✓:

Write a short inspiring poem about OpenAI, focusing on the recent DALL-E product launch (DALL-E is a text to image ML model) in the style of a {famous poet}



Reglas de oro para el uso de prompts

Evita imprecisiones y descripciones vagas

Less effective ✗:

The description for this product should be fairly short,
a few sentences only, and not too much more.

Better ✓:

Use a 3 to 5 sentence paragraph to describe this product.



Reglas de oro para el uso de prompts

Pon la instrucción al inicio del prompt y usa separadores (habitualmente """ o ###) para separar instrucción y contexto

Less effective ✗:

Summarize the text below as a bullet point list of the most important points.

{text input here}

Better ✓:

Summarize the text below as a bullet point list of the most important points.

Text: """

{text input here}

"""



Reglas de oro para el uso de prompts

Utiliza puntuación (ex., comillas) para clarificar el prompt

Less effective ✗:

```
Answer the given question using the given context. If no answer  
is found, return No answer.
```

```
Question: """{question}"""
```

```
Context: """{context}"""
```

Better ✓:

```
Answer the given question using the given context. If no answer  
is found, return "No answer".
```

```
Question: """{question}"""
```

```
Context: """{context}"""
```



Reglas de oro para el uso de prompts

En tareas específicas usa términos semilla para guiar la salida

Less effective ✗:

```
# Write a simple python function that
# 1. Ask me for a number in mile
# 2. It converts miles to kilometers
```

Better ✓:

```
# Write a simple python function that
# 1. Ask me for a number in mile
# 2. It converts miles to kilometers

import
```



Reglas de oro para el uso de prompts

Guía el formato de salida esperado

Less effective ✗:

Extract the entities mentioned in the text below. Extract the following 4 entity types: company names, people names, specific topics and themes.

Text: {text}

Show, and tell - the models respond better when shown specific format requirements. This also makes it easier to programmatically parse out multiple outputs reliably.

Better ✓:

Extract the important entities mentioned in the text below. First extract all company names, then extract all people names, then extract specific topics which fit the content and finally extract general overarching themes

Desired format:

Company names: <comma_separated_list_of_company_names>

People names: -||-

Specific topics: -||-

General themes: -||-

Text: {text}



Reglas de oro para el uso de prompts

Del zero-shot al few-shot

Zero-shot

Extract keywords from the below text.

Text: {text}

Keywords:

Few-shot - provide a couple of examples

Extract keywords from the corresponding texts below.

Text 1: Stripe provides APIs that web developers can use to integrate payment processing into their websites and mobile applications.

Keywords 1: Stripe, payment processing, APIs, web developers, websites, mobile applications

##

Text 2: OpenAI has trained cutting-edge language models that are very good at understanding and generating text. Our API provides access to these models and can be used to solve virtually any task that involves processing language.

Keywords 2: OpenAI, language models, text processing, API.

##

Text 3: {text}

Keywords 3:



Reglas de oro para el uso de prompts

Chain of Thought

Standard Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. 

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. 



Reglas de oro para el uso de prompts

Role prompting: asigna un rol al LLM que le ofrezca contexto sobre cómo actuar

Less effective ✗:

```
Please provide an analysis of the themes and  
narrative techniques in Jane Austen's "Pride and Prejudice."
```

Better ✓:

```
You are a literature professor who specializes in 19th-century  
English literature. Please provide an analysis of the themes and  
narrative techniques in Jane Austen's "Pride and Prejudice."
```



Fine-tuning LLMs

En caso de que el prompt no sea suficiente ⇒ Fine-tuning

⇒ Necesitas dataset de entrenamiento y capacidad computacional

- Herramientas de anotación NLP: [Argilla](#)

The image displays two screenshots of the Argilla annotation tool. The left screenshot shows a list of news articles with their text content and category labels (Business, Sci/Tech, Sports, World). The right screenshot shows a detailed view of a single record with a prompt, response, and feedback submission form.

⇒ TIP: usa OpenAI para generar la primer boceto de dataset de entrenamiento



Configuración de LLMs

Parámetros de optimización de LLMs

- Temperature:
 - Cuanto menor sea la temperatura, más deterministas serán los resultados
⇒ siempre se elige el siguiente token/palabra más probable.
 - Cuanto mayor la temperatura, más aleatoriedad y fomentar resultados más diversos o creativos
⇒ se aumenta la probabilidad de los otros posibles tokens/palabras
- Max tokens:
 - Limita el número máximo de tokens/palabras que puede generar el modelo



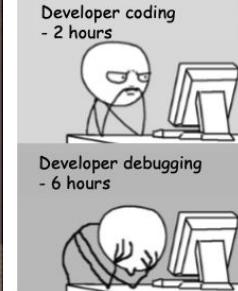
Conceptos LLM

Quieres aprender más:

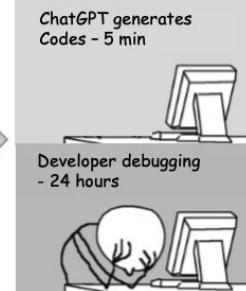
- [Prompt Engineering](#)
- [A Hackers' Guide to Language Models](#)
- [DeepLearning.AI courses](#)

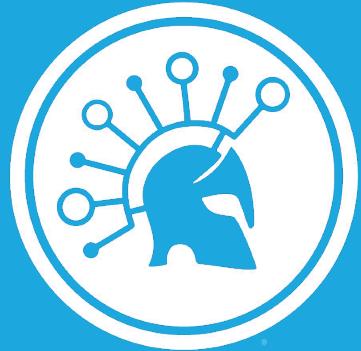


Days before OpenAI



Days after OpenAI





DataSpartan

Rétanos hoy

informacion@dataspartan.com

www.es.dataspartan.com