

Introducción al Procesamiento de Lenguaje Natural y Extracción de Información

Dr. Cristian Enrique Muñoz Villalobos

Agenda

- Que es el procesamiento de Lenguaje Natural?
- Aplicaciones y tareas de NLP
- Estrategias para la colecta de información
- Data Warehouse y Data Lakes
- Formatos abiertos: JSON, XML, HTML, etc.
- Web Crawling y APIs

"A language is not just words. It's a culture, a tradition, a unification of a community, a whole history that creates what a community is. It's all embodied in a language."

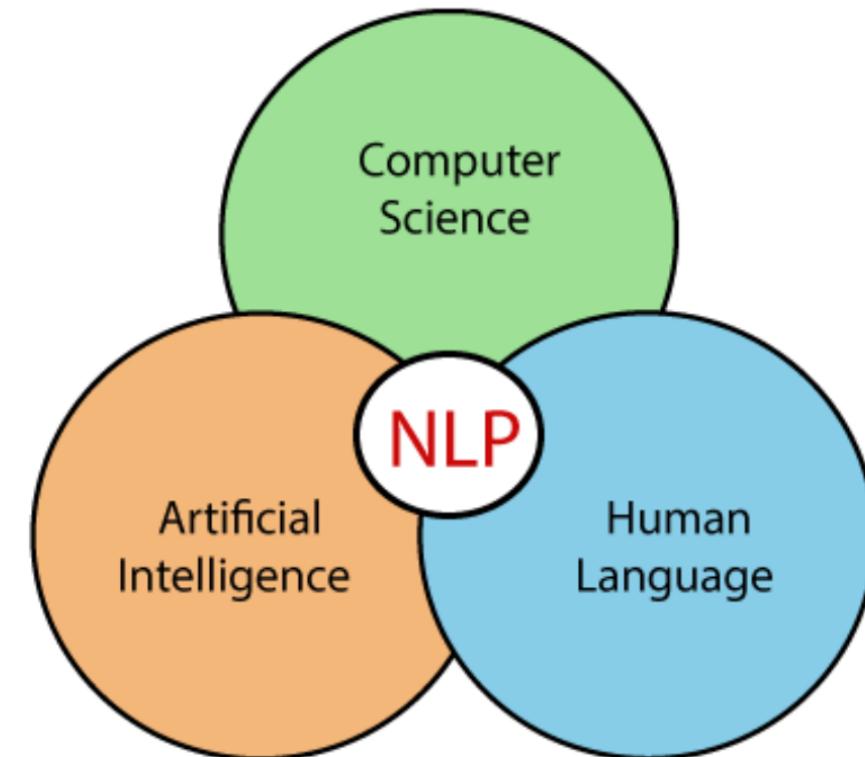
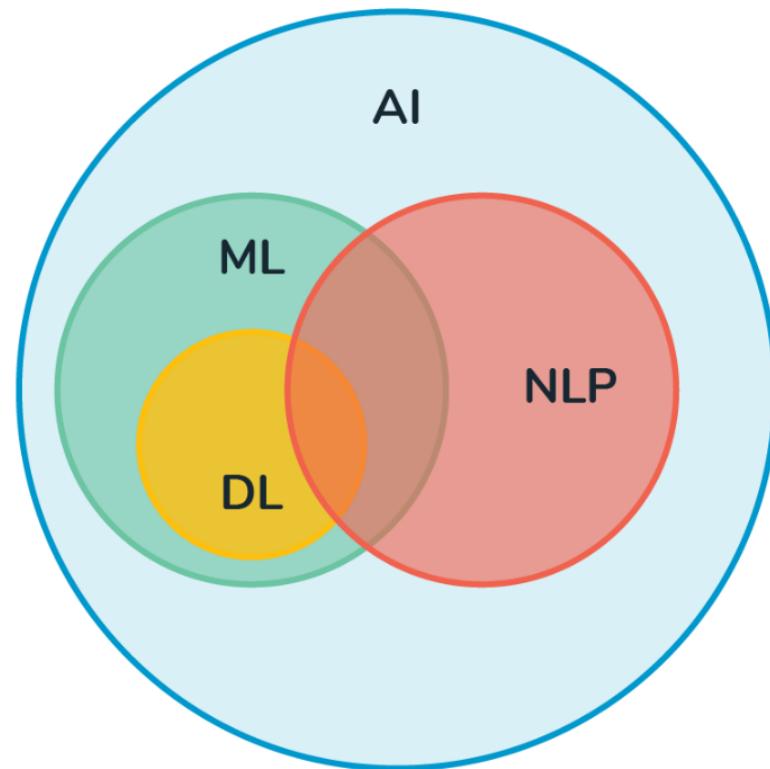
[Noam Chomsky]

Que es el Procesamiento de Lenguaje Natural?

Es un área de la inteligencia artificial que se especializa en el proceso computacional de entender y generar lenguaje humano.

Es una área de ciencia de la computación que utiliza métodos para analizar, modelar y entender el lenguaje humano.

Que es el Procesamiento de Lenguaje Natural?



Aplicaciones



Extracción de la información



*Clasificación de
Spam*



Análisis de Finanzas



*Recuperación de la
Información*



Motores de Búsqueda

*Extracción de Entidades
Legales*

Recuperación de la Información

Form 4562	Depreciation and Amortization (Including Information on Listed Property)	
Department of the Treasury Internal Revenue Service (IRS)	OMB No. 1545-0172	
	2020	Attachment Sequence No. 179
Name(s) shown on return	Business or activity to which this form relates	Identifying number
Part I Election To Expense Certain Property Under Section 179 Note: If you have any listed property, complete Part V before you complete Part I.		
1 Maximum amount (see instructions)	1	
2 Total cost of section 179 property placed in service (see instructions)	2	
3 Threshold cost of section 179 property before reduction in limitation (see instructions)	3	
4 Reduction in limitation. Subtract line 3 from line 2. If zero or less, enter -0-	4	
5 Dollar limitation for tax year. Subtract line 4 from line 1. If zero or less, enter -0-. If married filing separately, see instructions	5	
6 (a) Description of property	(b) Cost (business use only)	(c) Elected cost
7 Listed property. Enter the amount from line 29	7	
8 Total elected cost of section 179 property. Add amounts in column (c), lines 6 and 7	8	
9 Tentative deduction. Enter the smaller of line 5 or line 8	9	
10 Carryover of disallowed deduction from line 13 of your 2019 Form 4562	10	
11 Business income limitation. Enter the smaller of business income (not less than zero) or line 5. See instructions	11	
12 Section 179 expense deduction. Add lines 9 and 10, but don't enter more than line 11	12	
13 Carryover of disallowed deduction to 2021. Add lines 9 and 10, less line 12 ►	13	
Note: Don't use Part II or Part III below for listed property. Instead, use Part V. Part II Special Depreciation Allowance and Other Depreciation (Don't include listed property. See instructions.) 14 Special depreciation allowance for qualified property (other than listed property) placed in service during the tax year. See instructions. 15 Property subject to section 168(f)(1) election 16 Other depreciation (including ACRS)		
Part III MACRS Depreciation (Don't include listed property. See instructions.) Section A 17 MACRS deductions for assets placed in service in tax years beginning before 2020 18 If you are electing to group any assets placed in service during the tax year into one or more general asset accounts, check here □		
Section B—Assets Placed in Service During 2020 Tax Year Using the General Depreciation System (a) Classification of property (b) Month and year placed in service (c) Basis for depreciation (Business/investment use only—see instructions) (d) Recovery period (e) Convention (f) Method (g) Depreciation deduction 19a 3-year property b 5-year property c 7-year property d 10-year property		

Date/Time 03 Dec 2020

Payment Ref:
TID: 51572229

Total Amount: \$5.10

 Singapore Post Limited 10 Euros Road # Singapore Post Centre Singapore 408600	
Date/Time 03 Dec 2020 11:35 ID: S755202012030089 Location: Serangoon NEL Station No Description Amount 1. Postage Stamps Transaction No: S755203380015 Quantity: 1 Unit Price: \$0.60 Total GST: \$0.00 2. Postage Stamps Transaction No: S755203380015 Quantity: 3 Unit Price: \$1.50 Total GST: \$0.00	
Payment Ref: TID: 51572229 Total Amount: \$5.10 MID: Mode of Payment: VISA Auth Code: Card No: INV: 002482 ENT: Contactless Thank you for paying with VISA SAM Mobile enables you to make payments or manage your digital mail anytime, anywhere! Now available on Google Play and App Store	
 RECEIPT S.No. A- 3046195 KARMANAGA RANGA REDDY GOVERNMENT OF ANDHRA PRADESH eSeva 612748763 M.VITAYANAND PLOT NO. 39, DRAUGMOUR COLONY, SAROOR NADAK, SAROOR NADAK (SAROOR) 500007 Hyderabad Metro Water Works 0933244694 Bill for which Payment Made 27-09-2021 PAYMENT MODE: CASH Bank/Branch Name: WATER CESS Amount Paid (Rs.): 204.00 Amount in words (Rs.): Two Hundred Ninety Four Only Operator Code: Signature: 204	
 INVOICE - ONE SPEEDWAY BLVD. HOMESTEAD, FLORIDA 33036 305-247-1000 FEDERAL ID #05-077050 Invoice No. 230 Date: 03/03/2020 Order No.: Item No.: Event: GP 2020 Description: PLEASE PAY FROM THIS INVOICE Unit Price: TOTAL 1. THE FORGE CHRG DATED 3/21/20 \$460.00 \$460.00 1. THE FORGE CHRG DATED 3/21/20 \$13,642.00 \$13,642.00 1. EDWARD TILLOTSON CHRG DATED 3/21/20 \$700.00 \$700.00 Sub Total: \$14,742.00 Shipping & Handling: \$0.00 Taxes: \$0.00 Payment Method: CASH TOTAL: \$14,742.00	

Finanzas

DELL X Dataminr in partnership with ALL MARKET <1% NEWS/BLOG <1% CHATTER <1%

Streaming Apply Date Range

ALERTS

MN F Dell Says Special Committee Unanimously Determined That The Sale Of The Company Would Be The Best Alternative #Breaking 16:04

ALT Former HP CEO Carly Fiorina Says Dell Consolidation Possible via @besttechie @carlyforca @harveynashusa www.besttechie.com 16:58

MN Icahn builds Dell stake, complicating buyout: CNBC #MYTK4LIFE feeds.reuters.com 16:50

MKT Icahn Said Ready to Oppose Dell Deal dealbook.nytimes.com 16:45

MKT #Markets #News Icahn's Opposition to Dell LBO Could Doom Deal: Activist's 6% stake adds to ins... #Money #Barrons on.barrons.com 16:42

MN H-P, Southeastern Have Mulled Dell Bids via @WSJ 16:41

Summary 4:50AM - 5:09PM DELL ✓

MESSAGE GRAPH

400 200 0 16:00 03 AM 06 AM 09 AM 12 PM 03 PM \$14.40 \$14.20 \$14.00

MESSAGES

MKT: Dell Board Committee Insists Sale Was Best Outcome t.co/0ryezNrNz3 8:40 AM

MKT: BREAKING: Blackstone said to study Dell books in go-shop period \$BX 3:47 PM

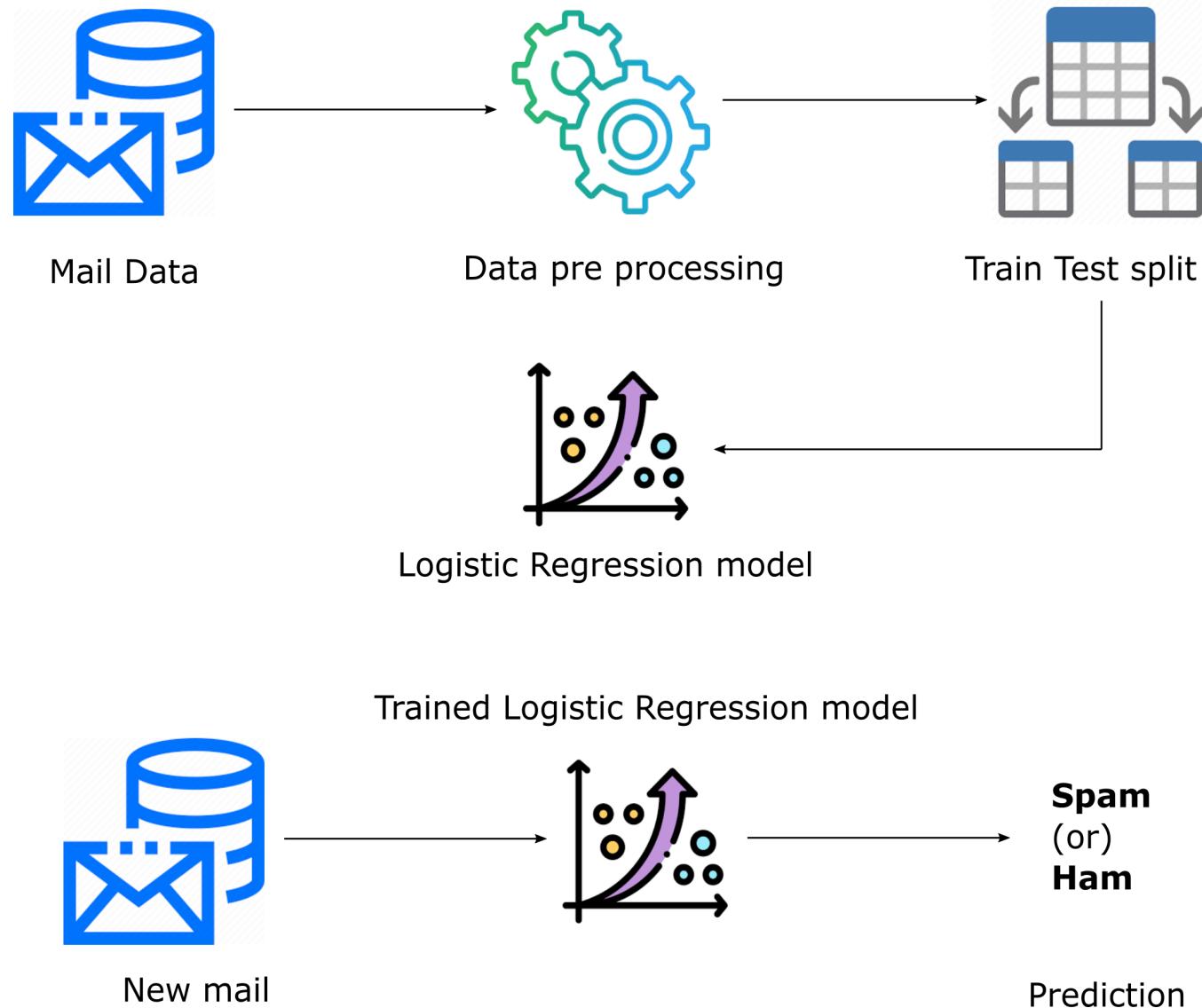
MN: Icahn builds up 6% stake in Dell on.ft.com/VHlwE2 #financial 4:43PM

MKT: Dell jumps on rumor; Cisco, H-P lead techs: t.co/HfNia13tEQ 3:08 PM

CLOUDS

DELL	HP	BUYOUT
CNBC	LENOVO	BOARD

Spam



Buscadores

Google

quién es el presidente de Perú

Todas Notícias Imagens Vídeos Maps Mais Ferramentas

Aproximadamente 114.000.000 resultados (0,53 segundos)

Peru / Presidente

Pedro Castillo

Pesquisas relacionadas

Keiko Fujimori Vlad... Aníbal Torres Alberto Fujimori Lilia Pare... Martín Vizca... Dina Bolu...

Feedback

<https://www.bbc.com> > mundo > ... ▾ Traduzir esta página

Quién es Pedro Castillo, el presidente electo de Perú, y en ...

20 de Jul. de 2021 — En una ajustadísima victoria, el candidato de izquierda Pedro Castillo, se impuso ante la candidata de derecha Keiko Fujimori en una de las ...

<https://www.bbc.com> > mundo > ... ▾ Traduzir esta página

Quién es Pedro Castillo, el maestro rural que llegó a la ... - BBC

20 de Jul. de 2021 — Con una imagen y un discurso que distan mucho de los de la tradicional élite política de Lima, Pedro Castillo logró finalmente ser ...



Chatbots

Chaty

Automatiza los chats de tu restaurante

Crea tu propio chatbot y dile adiós a tomar pedidos a mano

Probar demo

Beneficios

Facebook Messenger

WhatsApp Messenger

Hola, me gustaría realizar un pedido

Hola Daniel, bienvenido a "Mi restaurante"

¿Cómo puedo ayudarte? Escríbe una opción:

A. Realizar Pedido B. Horarios de Atención

Enviar mensaje

IBM Watson Assistant

Tareas de PLN

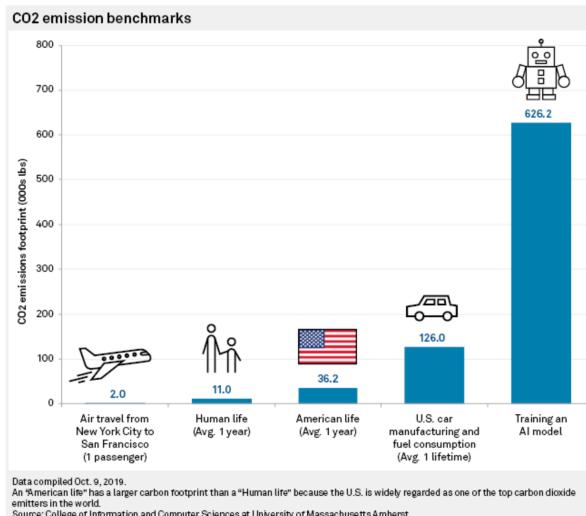
Modelamiento de Lenguaje	Clasificación de texto
Extracción de información	Recuperación de la información
Resúmenes de texto	principal del documento original
Agentes de Conversion	Pregunta-Respuesta
Maquinas de traducción	Moldeamiento de Tópicos

Modelos de Lenguaje

AI's large carbon footprint poses risks for big tech

The artificial intelligence industry has skyrocketed in recent years, powering technologies behind smart speakers and self-driving cars, but that growth is coming at a cost.

Researchers at the University of Massachusetts Amherst recently conducted a study assessing the energy consumption required to train several common large AI models. The study revealed that the training process can emit over 626,000 pounds of carbon dioxide, nearly 5x the lifetime emissions of an average car, or the equivalent of about 300 round-trip flights between New York and San Francisco.



VENTURE CAPITAL

A Golden Age For Natural Language

Louisa Xu Former Contributor

Dec 1, 2021, 11:00am EST

Listen to article 6 minutes

Every company talks to its customers with natural language. The last three years signaled the beginning of a golden age for natural language processing (NLP), one of the most useful and visible forms of machine learning (ML).

NLP is a branch of machine learning that endows computers with the human-like ability to understand text and spoken word. Thanks to the increase in computational power, the sheer volume of raw language data available on the Internet, and the popularization of deep learning, NLP has transformed from a hypothetical question posed in the 1950 Turing Test to an everyday reality.



ENTERPRISE TECH

Revolutionary NLP Model GPT-3 Poised To Redefine AI And Next Generation Of Startups

Hannah M. Mayer Contributor
Hannah is a former HBS & HSG researcher, and a Wiley-published author.

Jan 2, 2021, 03:09pm EST

Listen to article 9 minutes



LaMDA (Google Chatbot)

<https://www.washingtonpost.com/technology/2022/06/11/google-ai-lamda-blake-lemoine/>

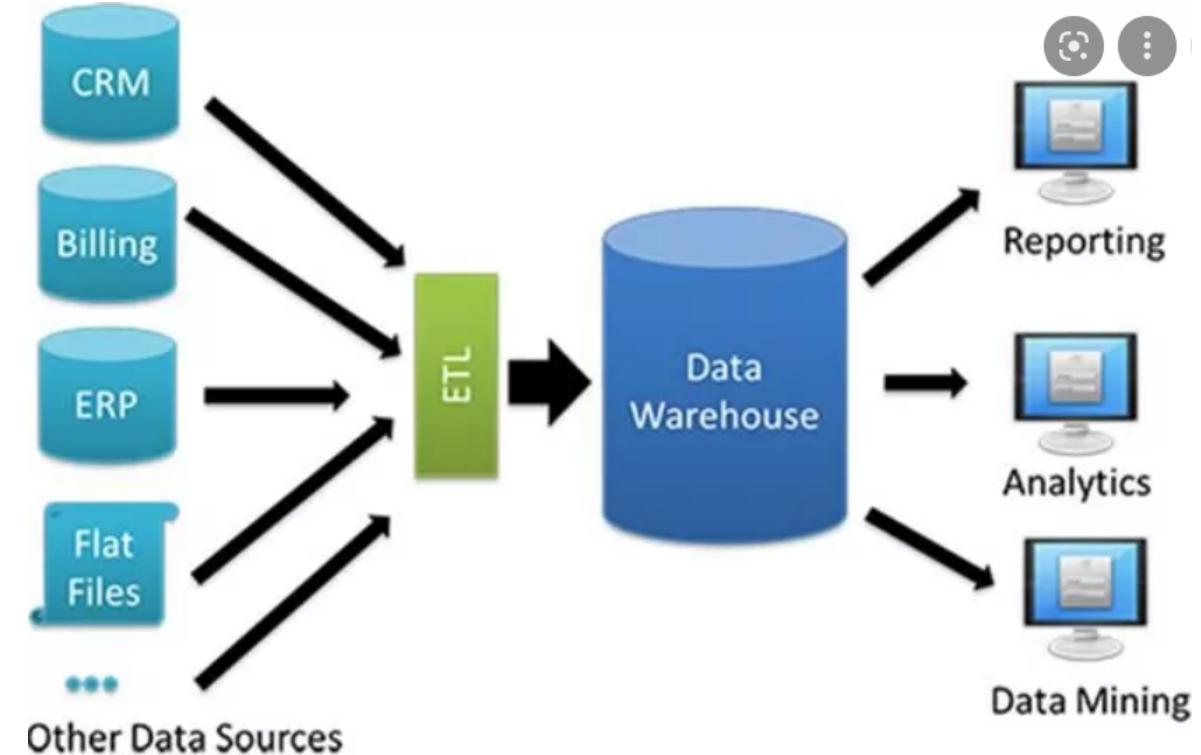
The screenshot shows a news article from the Washington Post. The title is "The Google engineer who thinks the company's AI has come to life". Below the title is a sub-headline: "AI ethicists warned Google not to impersonate humans. Now one of Google's own thinks there's a ghost in the machine." The author is listed as "By Natasha Tiku" on June 11, 2022, at 8:00 a.m. EDT. A large photo of Google engineer Blake Lemoine wearing a hoodie is on the left. To the right is a complex network diagram illustrating the AI's thought processes, with nodes and connections representing different thoughts and their relationships.

Estrategias para Colecta de Información

 Web Crawling	 Document Parsing
Recopilación automatizada de datos y su conversión en información estructurada para su análisis. e.g. datos de páginas web.	Implica la revisión de datos presentes en un documento y extraer información útil de él. e.g. PDF, xlsx, etc.

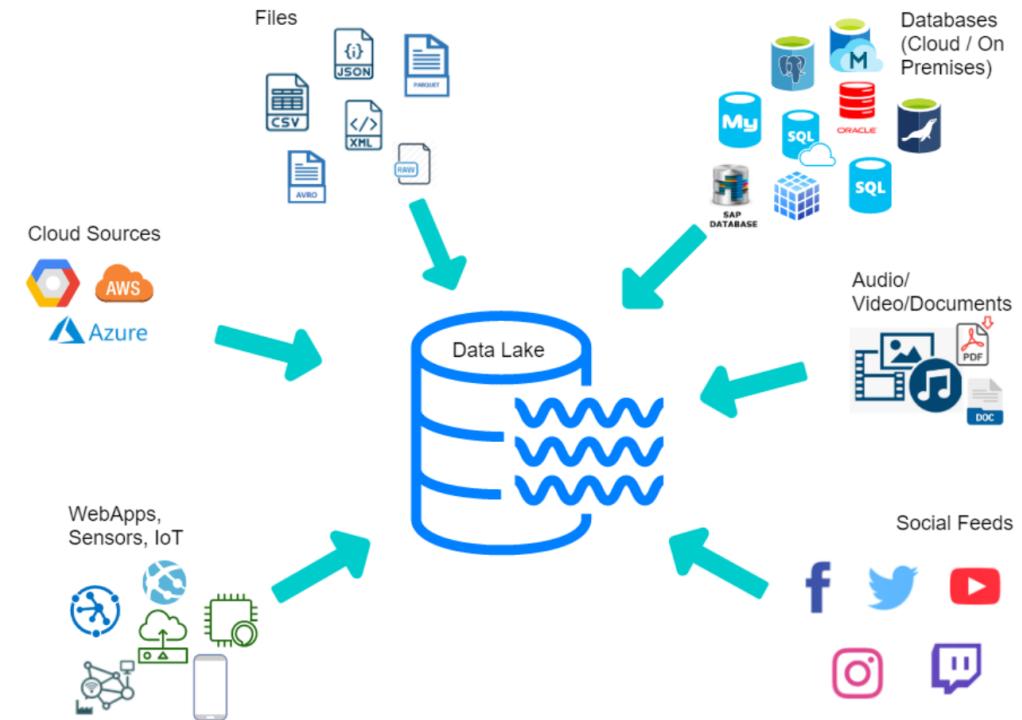
Data Warehouse y Data Lakes

Data Warehouse: es un banco de datos optimizado para analizar datos relacionales. La estructura de datos y los esquemas son definidos previamente para optimizar consultas SQL rápidas. Los datos son limpios, enriquecidos y transformados para que puedan actuar como la "fuente única verdadera" en que el usuario puede confiar.



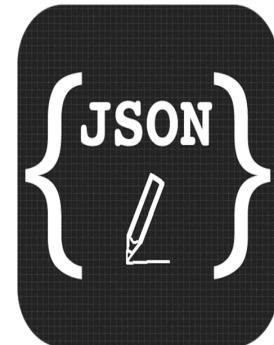
Data Warehouse y Data Lakes

Data Lake: son repositorios centralizados que permiten almacenar todos tus datos estructurados y no estructurados en cualquier escala. A principal ventaja de data lakes es la alta disponibilidad de los datos para aplicaciones de machine learning, análisis predictivas, descubrimiento de datos, etc.



Formatos Abiertos

Es un formato de archivo para almacenar datos digitales. Generalmente, mantenida por su organización de estándares y que puede ser utilizado e implementado por cualquier persona.



...

Formato JSON (Javascript Object Notation)

Relacionado al desarrollo de software, JSON es convertido en un patrón ampliamente utilizado para intercambio de información entre sistemas. Es rápido y simple de ser interpretado por los computadores y seres humanos. Google ha explorado ampliamente el uso de JSON en sus API. Muchos lenguajes de programación dan soporte a este formato: Python, Ruby, PHP, Javascript.

JSON Viewer: <http://jsonviewer.stack.hu/>

Official Site: <https://www.json.org/json-en.html>

```
{  
  "nombre": "Cristian Enrique Muñoz Villalobos",  
  "edad": 34,  
  "pais": "Perú",  
  "Meses": ["Enero", "Febrero", "Marzo"],  
  "Familia": {"Papá": "José", "Hermana": "Jannet"}  
}
```

Formato XML (Extensible Markup Language)

Archivo de texto considerado un padrón para intercambio de información entre sistemas y almacenamiento de datos. XML es un lenguaje de marcación definida por W3C.

XML Viewer: <https://jsonformatter.org/xml-viewer>

```
<pucp>
  <user id="1">
    <nOMBRE>Cristian Enrique Muñoz Villalobos</nOMBRE>
    <edad>34</edad>
    <pais>Perú</pais>
    <Meses>
      <Mes>Enero</Mes><Mes>Febrero</Mes>
    </Meses>
    <Familia>
      <papa>José</papa><hermana>Jannet</hermana>
    </Familia>
  </user>
</pucp>
```

Formato HTML (Hyper-Text Markup Language)

(1989) HTML esta en nuestra vidas todos los dias, mismo que no nos demos cuenta. Es un lenguaje de marcación. El formato utiliza **TAG** para estructurar a informa o que ser a exibida en la pagina web. (Obs: **CSS** por otro lado indica como el elemento estructural debera ser exibido: color/fuente/tam o/etc.)

- Tag HTML

```
<!-- Inicio y fin de pagina -->
<html></html>
```

- Tag HEAD

```
<head>
  <title>Introducción a NLP para Ciencias Sociales</title>
  <meta name="description" content="El mejor curso de NLP">
  <meta name="keywords" content="Reconocimiento de Entidades, Procesamiento de Lenguaje Natural, Extracción de Relaciones, Inteligencia Artificial">
</head>
```

- Tag Simples

```
<h1>Sobre el Curso</h1>
<p>Dirigido a Profesionales y estudiantes de carreras vinculadas a las Ciencias ...</p>
<p>Sumilla: Cantidades enormes de información relevante para las ciencias sociales son producidas diariamente ...</p>
<p>Página del curso: <a href="https://qlab.pucp.edu.pe/cursos/proc-lenguaje" target="_blank">NLP QLAB</a>.</p>
```

- Version 5

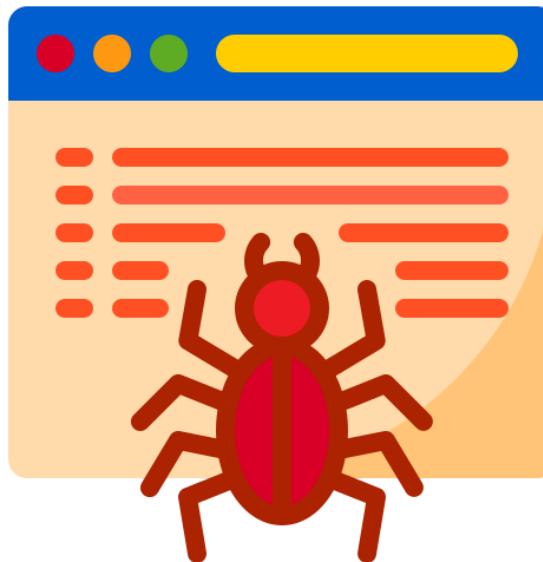
```
<!-- Elementos Semanticos -->
<header></header>
<footer></footer>
<article></article>
<section></section>

<!-- Elementos Gráficos -->
<svg></svg>
<canvas></canvas>

<!-- Elementos Multimedia -->
<audio></audio>
<video></video>
```

Como funciona Web Crawling?

También llamado bot o web spider es un algoritmo usado por los buscadores para encontrar, leer e indexar páginas web. Un ejemplo de crawler é o Googlebot de Google. Para aplicaciones específicas se recomienda el uso de bibliotecas open source como selenium o requests de python una flexibilidad mayor para recopilar información.



Redes Sociales

Las redes sociales son plataformas con alto interés para el análisis de información. Insight de productos, opiniones, tendencias son algunos ejemplos de información. La mayoría de redes sociales facilitan APIs que permite acceder de forma controlada, ordenada y eficiente a las informaciones de sus comunidades.



Seleinum (Python)

Esta biblioteca de permite de forma fácil acceder a todas las funcionalidades de Selenium WebDriver. Dispositivo que permite controlar um browser (Firefox, Chrome, IE) de forma remota.

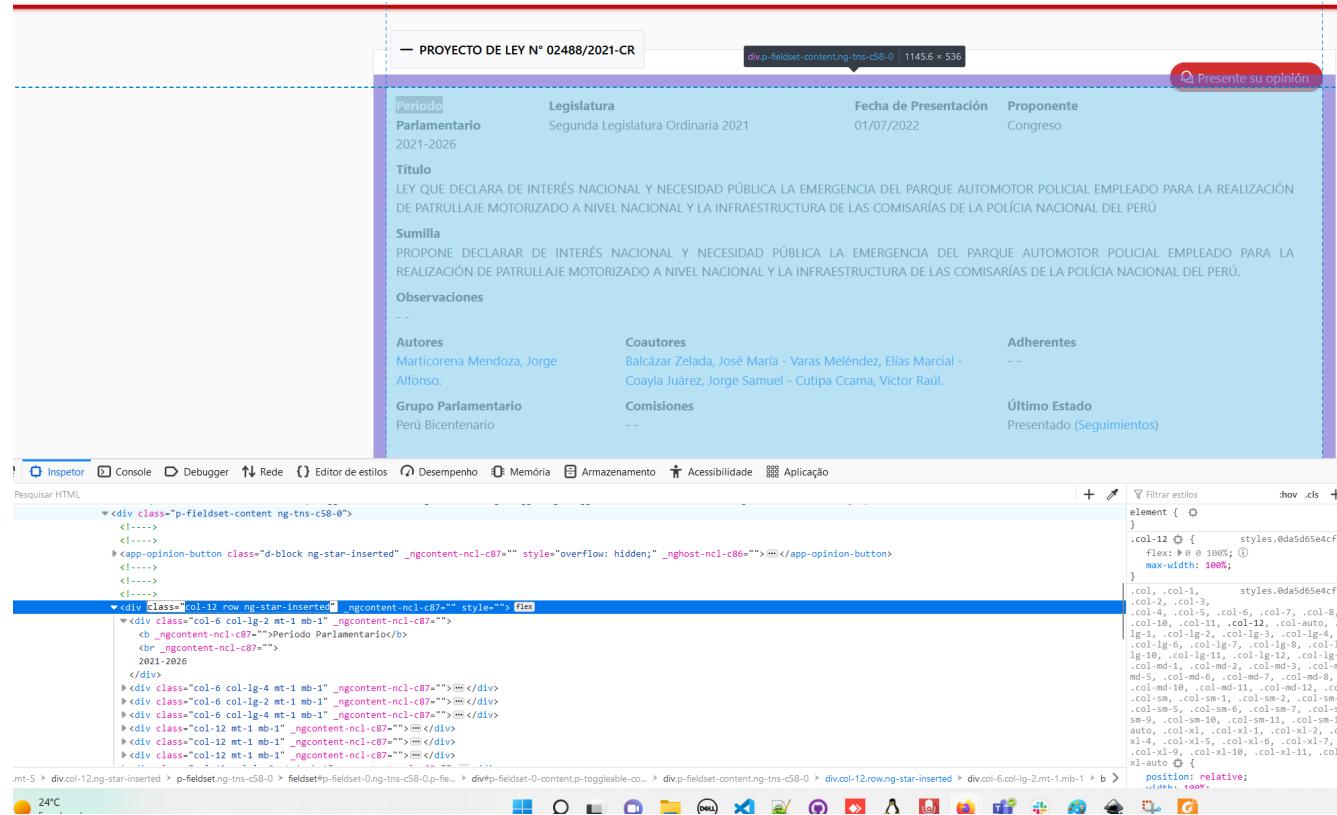
Site: <https://selenium-python.readthedocs.io/>

```
from selenium import webdriver
from selenium.webdriver.common.keys import Keys
from selenium.webdriver.common.by import By

driver = webdriver.Firefox()
driver.get("http://www.python.org")
elem = driver.find_element(By.NAME, "q")
driver.close()
```

Selenium (XPATH command)

```
//tagname[@attribute='value']  
elem = driver.find_element(By.XPATH, '//div[@class="col-12 row ng-star-inserted"]')
```



Tarea 1 :

Descarga de informaciones de proyectos de ley de la página del congreso de la república.