

**UNIVERSIDAD ESTATAL A DISTANCIA  
LABORATORIO DE INVESTIGACION E INNOVACIÓN TECNOLÓGICA  
PROYECTO PRO RED  
VICERRECTORÍA DE INVESTIGACIÓN**

**Proyecto:**  
FakeNews

**Informe y Reporte de Avances**  
*Manual de Text2graph*

**Estudiante Vinculado:**  
Brayan Rodríguez Delgado

**Coordinadores**  
  
Adriana Céspedes  
Andrés Segura Castillo

Noviembre 2020

## Tabla de Contenidos

1. Introducción:	3
2. Información General:	3
3. Estructura del proyecto:	3
4. Instalación:	3
4.1 Dependencias:	4
4.2 Diagrama básico del prototipo:	5
5. Extractor CrowdTangle (ct_extractor):	5
5.1. Objeto principal: get_cd_data(token):	5
5.2 Objetos ctpost, post y links:	6
5.3 Objeto posts:	7
5.4 Objeto lists:	8
5.5 Objeto list:	9
5.6 Objeto leaderboard:	9
5.7 Objeto search:	10
6. Tokenizador de textos (tokenizer_text):	10
6.1 tokenize_text:	11
7. Generador de grafos (data_graph):	11
7.1 Objeto data_graph:	12
7.2 Funciones internas, no muestran resultados para el usuario:	16
8. Capturador de texto html (html_text):	17
8.1. get_html_text:	17
8.2. search:	17
9. Mejoras a realizar:	18
10. Historial de Versiones:	18
Versión 0.0.2:	18

# Text2graph

## 1. Introducción:

Text2graph es una librería dirigida al análisis y visualización de grafos de conceptos, obtenidos de un texto, el cual puede ser directamente introducido por el usuario, mediante un enlace, leyendo un documento, una base de datos en .csv o bien capturando el texto por medio de CrowdTangle.

## 2. Información General:

**Nombre proyecto:** Text2graph

**Descripción:** Librería para la visualización y análisis de grafos generados mediante textos.

**Versión actual:** 0.0.2 (Prototipo)

**Autor o autores:** Brayan Rodríguez Delgado <bradrd2009jp@gmail.com>

**Institución:** Laboratorio de Investigación e Innovación Tecnológica LIIT-UNED.

**Coordinadores:** Adriana Céspedes y Andrés Segura.

## 3. Estructura del proyecto:

El proyecto cuenta con cuatro módulos independientes entre sí, aunque algunos módulos pueden alimentarse de la salida de los otros, estos módulos son: el generador de token (*tokenizer\_text*), el generador de grafos (*data\_graph*), un extractor de datos para CrowdTangle (*ct\_extractor*) y una librería para capturar texto mediante html (*html\_text*). La organización de la carpeta principal del proyecto presenta los siguientes archivos:

```
. /
| __ text2graph/
|   ____ ct_extractor.py
|   ____ data_graph.py
|   ____ html_text.py
|   ____ _init_.py
|   ____ tokenizer_text.py
|   __ AUTHORS.md
|   __ CrowdTangle_API.odt
|   __ install_model.py
|   __ LICENSE
|   __ MANIFEST.in
|   __ README.md
|   __ setup.cfg
|   __ setup.py
|   __ Text2Graph.odt
```

## 4. Instalación:

Para instalar la librería y sus dependencias, se requiere descomprimir los archivos principales, ingresamos a la carpeta donde se encuentre *text2graph\_master.zip*:

```
unzip text2graph_master.zip
```

o bien:

```
tar -xvzf text2graph_master.zip
```

Si se trabaja desde Windows, simplemente descomprimir el archivo proporcionado. Se ingresa en la carpeta descomprimida y se digita el siguiente comando:

```
python3 setup.py install --user
```

Si utiliza Windows puede utilizar el siguiente comando:

```
python setup.py install --user
```

Esta librería requiere instalar el modelo correspondiente de spacy, el cual se instala mediante los siguientes comandos en la terminal:

```
python3 install_model.py #Para linux
```

o bien

```
python install_model.py #En Windows
```

Se recomienda tener instalado python 3.6 o superior, además esta librería no se ha probado todavía en MacOS, por lo que no se puede brindar por el momento soporte para este sistema operativo.

#### 4.1 Dependencias:

Además de los archivos presentes en la carpeta principal del proyecto, también se instalarán las siguientes dependencias necesarias para el correcto funcionamiento de los módulos:

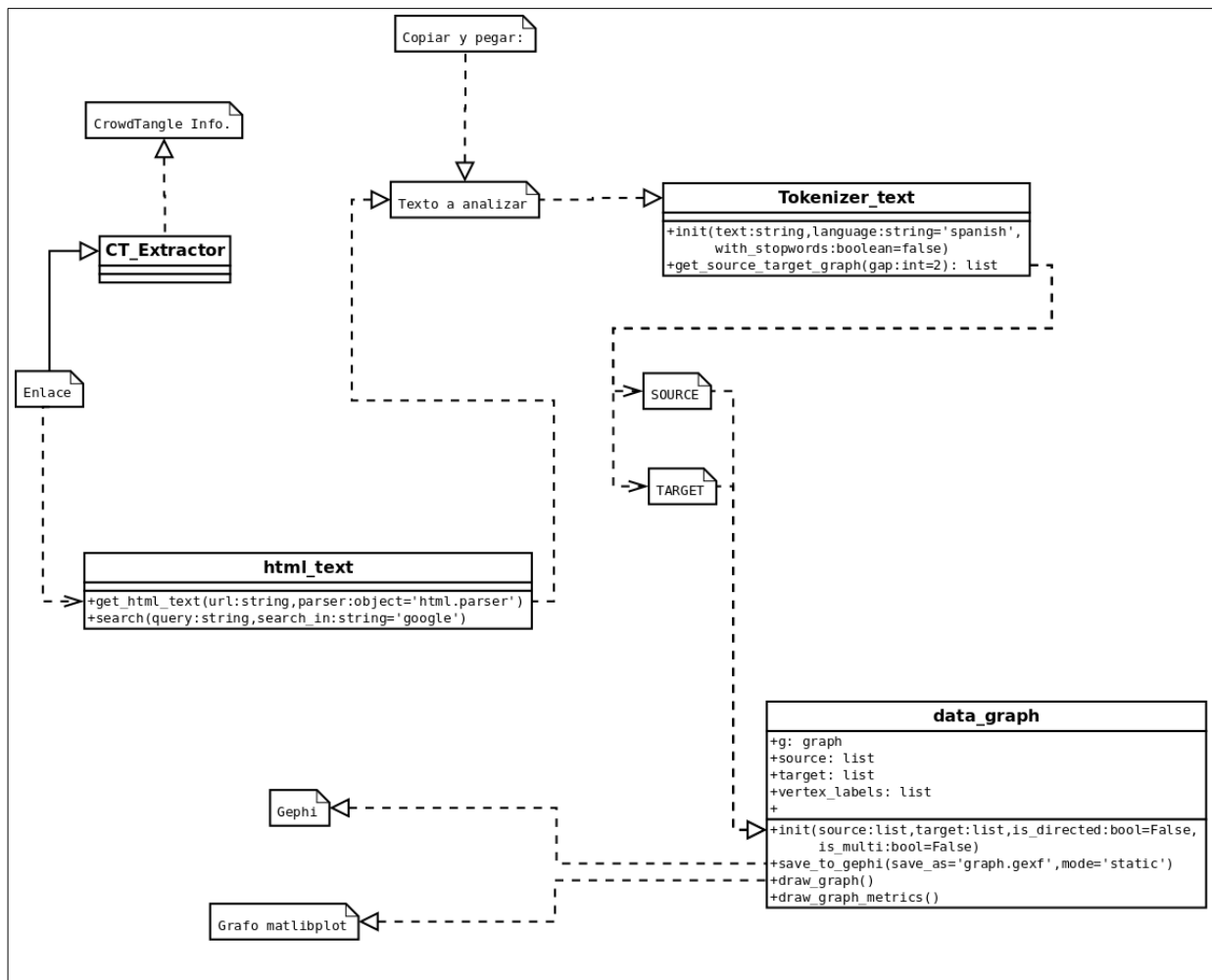
```
"networkx>=2.4",  
"pandas>=0.24.1",  
"pandas-datareader>=0.7.0",  
"plotly>=3.4.1",  
"scipy>=1.3.0",  
"numpy>=1.16.4",  
"pyreadr>=0.2.9",  
"spacy>=2.3.2",  
"tldextract>=2.2.3",  
"nltk>=3.4.0",  
"mpld3>=0.3.0",  
"beautifulsoup4>=4.9.1",  
"requests>=2.25.0",  
"tldextract>=2.2.3",  
"lxml>=4.3.1",  
"google>=3.0.0"
```

También es posible utilizar los modulos en forma independiente dentro de otros proyectos copiando los archivos de la carpeta text\_graph dentro del proyecto principal (requiere instalar las dependencias)

## 4.2 Diagrama básico del prototipo:

Text2graph cuenta con cuatro módulos principales o clases para realizar las funciones, separadas entre sí para permitir ser utilizadas como librería independiente cada uno de ellos, o en conjunto, o como parte de un proyecto mayor. Estas son:

- *ct\_extractor*: un extractor de datos de CrowdTangle
- *html\_text*: capturador de texto de una página web utilizando el enlace.
- *tokenizer\_text*: crea la lista de tokens de un texto y genera las listas de source y target para el grafo.
- *data\_graph*: que genera el grafo con la librería Networkx, y puede exportar a gephi.



**Figura 1.** Diagrama parcial del analizador de texto: text2graph

## 5. Extractor CrowdTangle (ct\_extractor):

Módulo que emula los comandos del API de CrowdTangle (Ver documento adjunto CrowdTangleAPI.odt), mediante el uso de clases y objetos que transforman las salidas *json* del API en datos de texto o csv según corresponda.

### 5.1. Objeto principal: get\_cd\_data(token)

get_cd_data(token)		
<b>Parámetros:</b>	token:	Crea un objeto de datos de CrowdTangle mediante el token dado

		por el API en el Dashboard principal.
<b>Salida:</b>	objeto:	<p>Devuelve un objeto CrowdTangle, con el cual se pueden realizar diferentes consultas, devolviendo los objetos correspondientes:</p> <p><i>ctpost</i>: por CrowdTangle post.  <i>post</i>: por Facebook post (id_ de Facebook).  <i>lists</i>: listas guardadas en el dashboard  <i>list</i>: devuelve la información de una lista mediante su número de identificación.  <i>links</i>: devuelve la información de CrowdTangle a través de un enlace o dirección, si existe una publicación asociada a Facebook a dicho enlace.  <i>posts</i>: devuelve la información de los posts almacenados dentro del leaderboard.  <i>leaderboard</i>: devuelve la información contendida en la tabla de clasificación principal del CrowdTangle.  <i>search</i>: devuelve la búsqueda mediante parámetros de CrowdTangle, requiere permiso de Facebook para poder realizarse este tipo de consulta, por lo que no se ha programado por completo.</p>

Ejemplo:

```
>>> import text2graph.ct_extractor as ce
>>> token = ...
>>> ct_dt = ce.get_ct_data(token)
```

## 5.2 Objetos *ctpost*, *post* y *links*:

<b>.ctpost(ctpost)</b>		
<b>Parámetros:</b>	ctpost:	Número de identificación dado por CrowdTangle al post.
<b>.post(fbpost, includeHistory=True)</b>		
<b>Parámetros:</b>	fbpost	Número de identificación de facebook [usuario]_[post]
	includeHistory	Si desea incluir historial de interacciones (True o False)
<b>.links(link, count=100, includeHistory=False, includeSummary=False, **kwargs)</b>		
<b>Parámetros:</b>	link	Enlace externo o interno de la publicación
	count	Número máximo de enlaces relacionados a mostrar, si los encuentra.
	includeHistory	Si se desea que se incluya el historial (True o False)
	includeSummary	Si se desea que se incluya el resumen (True o False)
	kwargs	startdate, endDate, sortBy, (consultar CrowdTangleAPI.odt)
<b>Funciones:</b>	.raw_dict():	devuelve el diccionario del <i>json</i> con la data original.
	.status():	devuelve el estado de la búsqueda: 200 si ha encontrado

		todos los datos.
	.result():	devuelve el diccionario con los resultados.
	.platform_id():	devuelve la plataforma del post.
	.date():	devuelve la fecha de publicación.
	.notes():	devuelve ntoaas si son incluidas en el post.
	.message():	devuelve el texto si es incluido en el post.
	.title():	devuelve el título si es incluido en el post.
	.ct_id():	devuelve el número de CrowdTangle.
	.link():	devuelve el link externo.
	.post_url():	devuelve el url dentro de Facebook.
	.domain():	devuelve el domino en facebook.
	.type():	devuelve el tipo de publicación (link, post)
	.media_type():	devuelve el tipo de media (video, foto, etc.)
	.media_url():	devuelve el url de media.
	.statisitcs():	devuelve el diccionario con las estadísticas de las interacciones.
	.statistics_df()	devuelve el dataframe con las estadísticas de las interacciones actuales y esperadas.
	.history():	devuelve el diccionario con el historial de las estadísticas.
	.history_df()	devuelve el dataframe con el historial de las estadísticas de las interacciones actuales y esperadas.
<b>Solo en link():</b>		
	.caption():	devuelve el subtítulo de la publicación.

Ejemplo:

```
>>> ctpost_numb = ...
>>> ct_dt.ctpost(ctpost_numb).statistics_df()
...
>>> fbpost_numb = ...
>>> ct_dt.post(fbpost_numb).history_df()
...
>>> link = 'http://. . .'
>>> domain = ct_dt.links(link).domain()
>>> print(domain)
...
```

### 5.3 Objeto posts:

.posts(count=10, includeHistory=False, includeSummary=False, **kwargs)		
<b>Parámetros:</b>	count:	Número de post a mostrar.

	includeHistory	Si incluye historial, admite dos valores True o False.
	includeSummary	Incluye resumen de los datos, (True o False).
	kwargs:	startdate, endDate, language, sortBy, types (consultar CrowdTangleAPI.odt)
<b>Funciones:</b>	.raw_dict():	devuelve el diccionario del <i>json</i> con la data original.
	.status():	devuelve el estado de la búsqueda: 200 si ha encontrado todos los datos.
	.result():	devuelve el diccionario con los resultados.
	.get_df():	devuelve el dataframe con la lista de post resumiento los siguientes datos:  <i>platformId</i> : tipo de plataforma. <i>date</i> : fecha. <i>update</i> : fecha actualización. <i>type</i> : tipo de publicación, foto, video, etc. <i>title</i> : título de la publicación. <i>caption</i> : subtítulo. <i>decription</i> : decripción. <i>message</i> : mensaje o texto. <i>link</i> : enlace externo. <i>postUrl</i> : url de facebook. <i>subscriberCount</i> : cantidad de subcriptores. <i>score</i> : puntaje en facebook.

Ejemplo:

```
>>> ctpost_numb = ...
>>> ct_dt.posts().get_df()
...
```

#### 5.4 Objeto lists:

.lists()		
<b>Parámetros:</b>	None:	
<b>Funciones:</b>	.raw_dict():	devuelve el diccionario del <i>json</i> con la data original.
	.status():	devuelve el estado de la búsqueda: 200 si fue exitosa.
	.result():	devuelve el raw_data de los resultados.
	.list_of_dict():	devuelve el diccionario con el resultado de búsqueda
	.list_df():	devuelve el dataframe con ID, Título y Tipo.
	.list_of_id():	devuelve la lista con los números de identificación de las listas.

Ejemplo:



```
>>> ct_dt.lists().list_df()
...
```

### 5.5 Objeto list:

.list(id_, count = 10, offset_options = 0)		
<b>Parámetros:</b>	id_:	Número de lista a mostrar.
	count	Número de renglones dentro de la lista a mostrar.
	offset_options	Opciones de impresión (Ver documento CrowdTangleAPI.odt).
<b>Funciones:</b>	.raw_dict():	devuelve el diccionario del <i>json</i> con la data original.
	.status():	devuelve el estado de la búsqueda.
	.result():	devuelve el diccionario de la búsqueda
	.accounts_df():	Devuelve el DataFrame con los siguiente atributos:  <i>id</i> : identificación. <i>name</i> : nombre de la cuenta. <i>handle</i> : manipulación. <i>profileImage</i> : imagen de perfil. <i>suscriberCount</i> : número de subscriptores. <i>url</i> : url de la cuenta. <i>platformId</i> : plataforma. <i>accountType</i> : tipo de cuenta. <i>pageAdminTopCountry</i> : lugar de administración. <i>verified</i> : si ha sido verificado.

### 5.6 Objeto leaderboard:

.leaderboard(count=50, **kwargs)		
<b>Parámetros:</b>	count:	Número de post a mostrar.
	kwargs:	startdate, endDate, orderBy, sortBy, (consultar CrowdTangleAPI.odt)
<b>Funciones:</b>	.raw_dict():	devuelve el diccionario del <i>json</i> con la data original.
	.status():	devuelve el estado de la búsqueda: 200 si ha encontrado todos los datos.
	.result():	devuelve el diccionario con los resultados.
	.get_df():	devuelve el dataframe con la lista de post resumiento los siguientes datos:  <i>account</i> : cuenta <i>summary</i> : resumen.

		<i>subscriberData</i> : datos del subcriptor. <i>id</i> : identificación. <i>name</i> : nombre de la cuenta. <i>handle</i> : manipulación. <i>suscriberCount</i> : número de subscriptores. <i>url</i> : url de la cuenta. <i>platformId</i> : plataforma. <i>pageAdminTopCountry</i> : lugar de administración. <i>verified</i> : si ha sido verificado. <i>likeCount</i> : total de likes <i>loveCount</i> : total de me encanta. <i>hahaCount</i> : total de me divierte. <i>wowCount</i> : total de me sorprende. <i>thankfulCount</i> : total de agradecido. <i>angryCount</i> : total de enojado. <i>sadCount</i> : total de me entristece. <i>shareCount</i> : total de compartido. <i>commentCount</i> : total de comentarios. <i>totalInteractionCount</i> : total de interacciones. <i>interactionRate</i> : tasa de interacciones.
--	--	--

Ejemplo:

```
>>> ct_dt.lists().leaderboard().get_df()
...
```

### 5.7 Objeto search:

.search(count=10, includeHistory=False, **kwargs)		
<b>Parámetros:</b>	count:	Número de post a mostrar.
	kwargs:	startdate, endDate, orderBy, sortBy, language, searchField, searchTerm (consultar CrowdTangleAPI.odt)
<b>Funciones:</b>	.raw_dict():	devuelve el diccionario del <i>json</i> con la data original.

No se ha programado más funciones debido a que se requiere permiso especial para utilizar la función search dentro del API de CrowdTangle.

## 6. Tokenizador de textos (tokenizer\_text):

Es módulo toma un texto, que puede ser introducido mediante un link, consumido desde CrowdTangle o copiado de alguna fuente externa y realiza una normalización del texto, eliminando tildes, mayúsculas, además de lematizar todos los radicales para formar el concepto base, y con ella realizar la lista de aristas para alimentar el módulo de grafos.

En esta versión, se está trabajando con una clase que devuelve el objeto de texto, una lista tokenizada y normalizada, que mediante los parámetros proporcionados por el usuario genera una lista de vértices de un grafo tanto de salida (*sources*) como de llegada (*targets*) que son utilizados por el generador de grafos (*data\_graph*).

## 6.1 tokenize\_text

tokenize_text(text, language='spanish', with_stopwords=False)		
<b>Parámetros:</b>	text:	Texto a ser tokenizado, puede provenir de diferentes fuentes.
	language	Configura el lenguaje a utilizar, por default se utiliza el español ( <i>spanish</i> ), para alimentar el modelo.
	with_stopwords	Mantiene las palabras que no agregan significado semántico como artículos, pronombres personales, preposiciones, etc.
<b>Funciones:</b>	.get_source_target_graph(gap=2):	
	<b>Descripción:</b>	Devuelve la lista de source y target para alimentar el modelo de redes, dada el número de saltos correspondientes: gap = 0: $p[0] \rightarrow p[1], \dots, p[n-1] \rightarrow p[n]$ gap = 1: $p[0] \rightarrow p[2], \dots, p[n-1-gap] \rightarrow p[n]$
	<b>Parámetros:</b>	Gap: o salto entre palabras
	<b>Salida:</b>	Source y target: listas con los nodos de salida y llegada respectivamente, cada renglón de las listas corresponde a una arista.

Ejemplo:

```
>>>import text2graph.tokenizer_text as tkt
>>>text = "En un lugar de la Mancha, de cuyo nombre no quiero acordarme, no ha mucho
tiempo que vivía un hidalgo de los de lanza en astillero, adarga antigua, rocín flaco y galgo
corredor. Una olla de algo más vaca que carnero, salpicón las más noches, duelos y
quebrantos los sábados, lantejas los viernes, algún palomino de añadidura los domingos,
consumían las tres partes de su hacienda."
>>>tokens = tkt.tokenize_text(text)
>>>gap = 2
>>>source, target = tokens.get_source_target_graph(gap=gap)
...
```

El resto de funciones del tokenizador de textos son utilizadas en los procesos internos y permiten limpiar el texto, filtrar *stopwords*, crear el radicalizador (*stemmer*), lematizar el conjunto de palabras.

## 7. Generador de grafos (data\_graph):

Este módulo genera un objeto tipo grafo, cuando recibe por parte del usuario la listas de nodos de salida o *source* y la de llegada o *target*, aunque es un módulo independiente de los anteriores, pueden utilizarse en conjunto para generar redes de conceptos y palabras a partir de un texto de determinado. Se utiliza como base la librería *networkx*, y posee un módulo que permite exportar las redes en formato “.gexf” para ser leído por aplicaciones como Gephi.

## 7.1 Objeto `data_graph`

<code>data_graph(source, target, is_directed=False, is_multi=False)</code>		
<b>Parámetros:</b>	<code>source</code>	Lista con los nodos de salida.
	<code>target</code>	Lista con los nodos de llegada.
	<code>is_directed</code>	Si desea crear un grafo dirigido (True), o bien un grafo no es dirigido (False).
	<code>is_multi</code>	Si se desea crear un multigrafo (True), o bien un grafo simple (False).
<b>Atributos:</b>	<code>.g</code>	devuelve el grafo interno generado para ser trabajado por la librería networkx.
	<code>.source</code>	devuelve la lista de nodos de salida.
	<code>.target</code>	devuelve la lista de nodos de llegada.
	<code>.vertex_labels</code>	devuelve la lista con las etiquetas para los nodos.
<b>Funciones:</b>	<code>.get_node_frequency_df(save=False, save_as='node_freq.csv', index=False)</code>	
	<b>Descripción:</b>	Genera el dataframe con la frecuencia de los nodos tanto en source como en target. Se puede manipular mediante la librería pandas.
	<b>Parámetros:</b>	<code>save</code> : si desea que sea salvado a un archivo csv. <code>save_as</code> : Nombre para ser salvado por default se tiene "node_freq.csv" <code>index</code> : omite el índice en el dataframe guardado.
	<b>Salida:</b>	Source y target: listas con los nodos de salida y llegada respectivamente, cada renglón de las listas corresponde a una arista.
	<code>.plot_node_frequency(tail_number=25, save=False, save_as='node_freq.png', html=False, ascending=True, html_name='node_freq.html', figw = 10, figh = 10, color='g')</code>	
	<b>Descripción:</b>	Crea un gráfico de barras horizontales con la frecuencia de los nodos.
	<b>Parámetros:</b>	<code>tail_number</code> : número de elementos a mostrar. <code>save</code> : si desea que sea salvado (True), o no (False). <code>save_as</code> : nombre para guardar el archivo. <code>html</code> : si desea exportar a html (True), o no (False). <code>html_name</code> : nombre para guardar el html. <code>ascending</code> : de menor a mayor. <code>figw</code> : ancho de la figura. <code>figh</code> : altura de la figura. <code>color</code> : color de las barras.
	<b>Salida:</b>	imagen, html con el gráfico de barras horizontales de frecuencia.
	<code>.self_loop_nodes_number()</code>	
	<b>Descripción:</b>	Devuelve el número de nodos con conexión a sí mismos.

	<b>Parámetros:</b>	
	<b>Salida:</b>	Entero.
	<code>.adj_matrix_dataframe()</code>	
	<b>Descripción:</b>	Devuelve el dataframe con la matriz de adyacencia, incluye las etiquetas de los nodos.
	<b>Parámetros:</b>	
	<b>Salida:</b>	Dataframe manipulable con librería pandas.
	<code>.adj_matrix_csv(save_as='adjacency.csv', index=False)</code>	
	<b>Descripción:</b>	Almacena en disco el dataframe de la matriz de adyacencia.
	<b>Parámetros:</b>	<i>save_as</i> : nombre para salvar la matriz de adyacencia. <i>index</i> : False si no quiere que aparezca el índice en el dataframe.
	<b>Salida:</b>	.csv
	<code>.number_of_vertex()</code>	
	<b>Descripción:</b>	Devuelve el número de vértices o nodos.
	<b>Parámetros:</b>	
	<b>Salida:</b>	Entero.
	<code>.number_of_edge()</code>	
	<b>Descripción:</b>	Devuelve el número de aristas.
	<b>Parámetros:</b>	
	<b>Salida:</b>	Entero.
	<code>.shortest_distance_matrix()</code>	
	<b>Descripción:</b>	Devuelve la matriz con la distancia más corta entre nodos.
	<b>Parámetros:</b>	
	<b>Salida:</b>	Array, manipulable con la libería numpy.
	<code>.shortest_distance_dataframe()</code>	
	<b>Descripción:</b>	Devuelve la matriz con la distancia en formato dataframe.
	<b>Parámetros:</b>	
	<b>Salida:</b>	Data frame, manipulable con la libería pandas.
	<code>.metrics_df(metrics='all')</code>	
	<b>Descripción</b>	Devuelve el dataframe con todas la métricas indicadas por el usuario.
	<b>Parámetros</b>	<i>metrics</i> : que puede ser 'all' y devuelve todas las métricas programadas, o bien una lista con los siguientes nombres: <i>metrics</i> = ['degree', 'eccentricity', 'pagerank', 'eigenvector', 'harmonic', 'closeness', 'betweenness']
	<b>Salida:</b>	Dataframe manipulable con librería pandas.

	<code>.metrics_csv( save_as='metrics.csv', metrics='all', index=False)</code>	
	<b>Descripción</b>	Devuelve el dataframe y lo guarda como .csv.
	<b>Parámetros</b>	<i>metrics</i> : que puede ser 'all' y devuelve todas las métricas programadas, o bien una lista con los siguientes nombres: <i>metrics</i> = ['degree', 'eccentricity', 'pagerank', 'eigenvector', 'harmonic', 'closeness', 'betweenness'] <i>save_as</i> : nombre para guardar el dataframe. <i>index</i> : ignora el índice del dataframe en el documento .csv.
	<b>Salida</b>	.csv
	<code>.plot_node_metric(metric, tail_number=25, ascending=True, save=False, save_as='metric.png', html=False, html_name='metric.html', figw = 10, figh = 10, color='g')</code>	
	<b>Descripción</b>	Muestra el gráfico de los nodos según la métrica elegida por el usuario.
	<b>Parámetros</b>	<i>metric</i> : métrica elegida por el usuario, entre 'degree', 'indegree', 'outdegree', 'eccentricity', 'pagerank', 'eigenvector', 'harmonic', 'closeness', 'betweenness' <i>tail_number</i> : número de elementos a mostrar. <i>save</i> : si desea que sea salvado (True), o no (False). <i>save_as</i> : nombre para guardar el archivo. <i>html</i> : si desea exportar a html (True), o no (False). <i>html_name</i> : nombre para guardar el html. <i>ascending</i> : de menor a mayor. <i>figw</i> : ancho de la figura. <i>figh</i> : altura de la figura. <i>color</i> : color de las barras.
	<b>Salida</b>	Figura en html o png, o mostrar en pantalla.
	<code>.save_to_gephi(save_as='graph.gexf', description="", mode='static')</code>	
	<b>Descripción</b>	Salva el grafo a formato .gexf para ser leído por gephi.
	<b>Parámetros</b>	<i>save_as</i> : nombre del archivo. <i>description</i> : descripción dada por el usuario. <i>mode</i> : 'static' o 'dynamic'
	<b>Salida</b>	Archivo .gexf
	<code>.set_nx_layout( layout='random', **kwargs)</code>	
	<b>Descripción</b>	Configura el layout para mostrar el grafo.
	<b>Parámetros</b>	<i>layout</i> : que puede ser shell, spring, circular, kamada_kawai, bipartite, spectral, spiral, o planar. <i>kwargs</i> : argumentos específicos a cada layout, consultar documentación de Networkx.
	<b>Salida</b>	pos
	<code>.draw_graph(g = None, html=False, save=False, html_name = 'graph1.html', save_as='graph1.png', **kwargs)</code>	
	<b>Descripción:</b>	dibuja el grafo con la librería matplotlib, admite

		diferentes parámetros.
	<b>Parámetros</b>	<i>g</i> : grafo <i>html</i> : si desea salvar a html (True) o no (False) <i>save</i> : si desea salvar a png (True) o no (False) <i>html_name</i> : nombre para guardar html. <i>save_as</i> : nombre para guardar a png. <i>**kwargs</i> : otros parámetros
	<b>Salida</b>	grafo
	<pre>.draw_graph_metrics(     html=False,     save=False,     html_name='graph1.html',     save_as='graph1.png',     metrics='pagerank',     size_factor = 10000,     with_labels = False,     with_node_label = True,     with_values = False,     first_n_values = 5,     normal_node_size = True,     cmap = 'YlOrRd',     **kwargs):</pre>	
	<b>Descripción:</b>	Dibuja y colorea el grafo según una métrica específica
	<b>Parámetros</b>	<i>html</i> : si se salva a html, True, si no False. <i>save</i> : si se salva a png, True, si no False. <i>html_name</i> : nombre con que se guarda el html. <i>save_as</i> : nombre con que se guarda el png. <i>metrics</i> : tipo de métrica, 'pagerank', etc. <i>size_factor</i> : factor de multiplicación por la métrica <i>with_labels</i> : etiquetas en el grafo <i>with_node_label</i> : nodos con etiqueta <i>with_values</i> : nodos con valores de la métrica <i>first_n_values</i> : etiqueta para los primeros n valores. <i>normal_node_size</i> : tamaño normal del nodo (True) <i>cmap</i> : tipo de mapeo de color, consultar documentación de matplotlib.
	<b>Salida</b>	grafo
	.remove_node_by_label(label)	
	<b>Descripción:</b>	Remueve un nodo conociendo su etiqueta.
	<b>Parámetros</b>	<i>label</i> : la etiqueta del nodo a ser eliminado.
	<b>Salida</b>	
	.reset()	
	<b>Descripción:</b>	devuelve el grafo a su condición inicial.
	<b>Parámetros</b>	
	<b>Salida</b>	

	<code>.add_label_attrib()</code>	
	<b>Descripción:</b>	agrega el atributo de etiquetas a los nodos, cuando se trabaja con el exportador <code>.gexf</code> de <code>networkx</code> , las etiquetas deben pasar mediante atributos de nodos.
	<b>Parámetros</b>	
	<b>Salida</b>	
	<code>set_node_attributes(attr_dict, name)</code>	
	<b>Descripción:</b>	agrega un atributo a un nodo definido por el usuario, útil cuando se exporta el <code>.gexf</code> , ya interno o bien el de <code>networkx</code> .
	<b>Parámetros</b>	<i>attr_dict</i> : diccionario basado en la identificación de nodos dada por <code>.graph().nodes</code> o <code>.g.nodes</code> y el atributo correspondiente. <i>name</i> : nombre que desea que tenga el atributo.
	<b>Salida</b>	
	<code>.write_gexf(save_as='g_graph.gexf')</code>	
	<b>Descripción</b>	llamada al módulo de <code>Networkx</code> para exportar a <code>.gexf</code> , requiere que los nodos tengan propiedades para exportarlas, en especial las etiquetas.
	<b>Parámetros</b>	<i>save_as</i> : nombre con el que el archivo aparecerá.
	<b>Salida</b>	

Ejemplo:

```
>>> import text2graph.data_graph as tg
>>> text_graph = tg(source, target)
>>> text_graph.plot_node_frequency()
>>> text_graph.save_to_gephi()
>>> text_graph.set_nx_layout(layout='spring')
>>> text_graph.draw_graph_metrics(metrics='betwenness')
...
```

## 7.2 Funciones internas, no muestran resultados para el usuario:

<code>.create_graph()</code>	la función interna <code>.create_graph()</code> genera el grafo con el <code>source</code> y <code>target</code> brindados por el usuario, se almacena en la variable <code>g</code> .
<code>.graph()</code>	devuelve el grafo almacenado en <code>g</code> .
<code>.labels()</code>	devuelve el diccionario con las etiquetas de los diferentes nodos.
<code>.get_adj_matrix()</code>	devuelve la matriz de adjacencia de los nodos del grafo, no incluye etiquetas.
<code>.empty_dic()</code>	crea un diccionario vacío.
<code>.indegree()</code>	devuelve el diccionario con el <i>indegree</i> de los nodos, solo se calcula en grafos dirigidos.
<code>.outdegree()</code>	devuelve el diccionario con el <i>outdegree</i> de los nodos, solo se calcula en grafos dirigidos.



.eccentricity()	devuelve el diccionario con la <i>eccentricity</i> de los nodos, solo se calcula en grafos dirigidos, solo es posible calcularlo en grafos conectados.
.harmonic centrality()	devuelve el diccionario con la centralidad armónica, si no es posible calcularla devuelve un diccionario vacío.
.harmonic_closeness()	devuelve el diccionario con la centralidad de cercanía, si no es posible calcularla devuelve un diccionario vacío.
.eigenvector centrality()	devuelve el diccionario con la centralidad por eigenvector, con alpha 0.85 y un épsilon de 1e-3.
.betweenness()	devuelve el diccionario con la centralidad por intermediación normalizada.
.return_extension()	devuelve la extensión de un archivo
.default_layout()	configura el layout del grafo.
.graph_pos()	devuelve la configuración del layout.
.remove_node(node)	remueve el nodo según el número de nodo.
.list_nodes_attributes()	lista con los nombres de los atributos definidos por el usuario.
.get_dict_node_attrib()	diccionario con clave nombre de los atributos y el atributo correspondiente.

## 8. Capturador de texto html (html\_text):

Este módulo, permite mediante un link o enlace capturar todo el texto del html para poder alimentar el tokenizador en forma automatizada.

### 8.1. get\_html\_text:

get_html_text(url, parser='html.parser')		
<b>Parámetros:</b>	url:	dirección de html del cual se va a extraer el texto de la noticia.
	parser:	configura el extractor de etiquetas html de <i>BeautifulSoup</i> , puede ser <i>html.parser</i> , <i>lxml</i> , <i>xml</i> o <i>html5lib</i> , para mayor detalle consultar documentación de BeautifulSoup en: ( <a href="https://www.crummy.com/software/BeautifulSoup/bs4/doc/">https://www.crummy.com/software/BeautifulSoup/bs4/doc/</a> )
<b>Salida:</b>	output:	Texto del html contenido entre los tags <p></p>

Ejemplo:

```
>>> import text2graph.html_text as html
>>> text = html.get_html_text('https://en.wikipedia.org/wiki/Lorem_ipsum')
>>> print(text)
...
```

### 8.2. search:

search(query, search_in='google', timeout=None, **kwargs)		
<b>Parámetros:</b>	query:	Palabras claves de búsqueda.
	search_in:	Selecciona la máquina de búsqueda 'google' o 'yahoo'
	timeout:	Indica el tiempo de espera de la búsqueda.
<b>Salida:</b>	df:	Devuelve el DataFrame con los datos de la búsqueda, con título del html y el url de la búsqueda.

Ejemplo:

```
>>> import text2graph.html_text as html
>>> df = html.search('covid', search_in='google')
>>> print(df)
>>> text = html.get_html_text(df['link'][24])
...
```

## 9.Mejoras a realizar:

- En el futuro se puede programar el módulo de search del ct\_extractor, para tener las funcionalidades completas en el módulo.
- El lematizador requiere de mejorar en algunos aspectos, posiblemente integrarse o bien utilizar los modelos de los proyectos paralelos realizados por el equipo de #FakeNews.
- El módulo de métricas de *data\_graph*, puede recibir parámetros específicos para cada una de las métricas.
- El código del tokenizador de texto puede mejorarse, en especial en lo que se refiere al uso del español como lenguaje predeterminado. Es necesario encontrar una manera óptima para el manejo de las excepciones del idioma español. Existen problemas con la lematización correcta de los verbos irregulares y la clasificación de ciertos sustantivos, la radicación no es suficiente.
- Falta por decidir se es mejor trabajar los grafos con radicales o seguir trabajando con lematización.
- Actualmente el exportador de .gexf lo hace a la version 1.2, para asegurar el máximo de compatibilidad con gephi, pero se espera que en versiones posteriores se exporte al actual estándar 1.3.
- En el capturador de texto de html, falta mejorar el filtrado del texto, puede capturar algunos textos como etiquetas de comentarios o similares que puede no estar en el cuerpo principal.

## 10.Historial de Versiones:

Cambios principales entre versiones:

*Versión 0.0.2:*

- Se independiza el generador de grafos (*data\_graph*) del tokenizador (*tokenizer\_text*).
- Se agrega extractor de texto de html (*html\_text*).
- Se agrega un buscador de términos para las máquinas de búsqueda Google y Yahoo, devuelve los *url* y títulos de las búsquedas.

## Documentación consultada:

Gephi.org. (2007). **Gexf file Format**. Documentación y formato de archivos gexf. Consultado en:

<https://gephi.org/gexf/format/>

Richarson, L. (2020) **BS4: BeautifulSoup**. Documentación. Consultado en:  
<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>