



INSTITUTO POLITÉCNICO NACIONAL

---

ESCUELA SUPERIOR DE COMPUTO

*TRABAJO TERMINAL*

**RECOLECTOR Y CLASIFICADOR DE  
NOTICIAS**

**2018-B013**

**PRESENTAN:**

CARLOS ANDRES HERNANDEZ GOMEZ  
LUIS DANIEL MEZA MARTINEZ

**DIRECTORES:**

**M. en C. JOEL OMAR JUÁRES GAMBINO**  
**Dra. CONSULO VARINIA GARCIA MENDOZA**

**CIUDAD DE MÉXICO**

---

---

# Índice general

<b>1. Introducción</b>	<b>1</b>
1.1. Problemática . . . . .	2
1.2. Solución Propuesta . . . . .	2
1.3. Objetivo . . . . .	2
1.4. Objetivos Específicos . . . . .	2
1.5. Estructura del Documento . . . . .	3
<b>2. Conclusiones</b>	<b>5</b>
<b>3. Bibliografía</b>	<b>7</b>
3.1. Referencias . . . . .	7



# Capítulo 1

## Introducción

El objetivo de estudio de este trabajo es clasificar noticias implementando procesamiento de lenguaje natural, algoritmos de aprendizaje automático y el manejo tecnologías web (Crawler). El artículo periodístico es la información de un hecho ocurrido durante un lapso determinado, permite conocer el estado económico de un país, avances en la ciencia, desastres naturales, nivel de inseguridad, dichos acontecimientos independientemente del tema, día y lugar en el cual ocurrieron, tienen un impacto en la sociedad. Las características principales de este tipo de información es depender de un medio de comunicación como la televisión, redes sociales, diarios, blogs, los cuales crean expectativas en los sectores económicos (Educación, industria, turismo, etc.), modifica los planes de inversión de las empresas o naciones, siendo así de suma importancia su distribución de una forma eficaz. Cabe mencionar que el medio más utilizado es la internet.

El uso de las páginas web está en incremento, permitiendo consultar noticias de distintos sitios como los periódicos electrónicos; su información al igual que un diario tradicional se encuentra dividida en secciones para facilitar la consulta, sin embargo, la clasificación suele variar en cada compañía de prensa, incluso con el mismo contenido, un problema mayor se encuentra en los sitios independientes, los cuales no cuentan con una segmentación particular, haciendo difícil realizar una búsqueda eficaz.

Se han seleccionado los diarios más utilizados en México [1], con una buena segmentación en su contenido y se ha homogenizado las secciones en común, para obtener los datos necesarios (Noticias clasificadas) y realizar el entrenamiento del algoritmo de clasificación.

## 1.1. Problemática

Hoy en día existen distintas maneras de informarse acerca de los acontecimientos más recientes, por ejemplo, en la televisión, blogs, redes sociales, foros, diarios, etc. Esto ha provocado que la información se encuentre más dispersa y se tenga que acceder a muchos recursos para recopilarla. Esta situación implica un gran esfuerzo y tiempo. Para ayudar en este problema existen herramientas que hacen la búsqueda de noticias de interés para el usuario en forma automática. Sin embargo, dichas herramientas dependen de que los sitios a consultar cuenten con una etiquetación correcta y homogénea de la información.

Según El Economista [8] el sitio web “Animal Político” ocupa el lugar número cuatro en el ranking de medios nativos digitales y clasifica sus noticias de una manera poco habitual para los lectores, como la sección “El sabueso”, “El plumaje”, “Hablemos de . . .”, entre otras, lo que hace complicado obtener los artículos para los métodos tradicionales de recopilación que se basan sólo en las etiquetas que identifican cada sección y no el contenido de las noticias. Debido a lo anterior se propone crear un recolector de noticias el cual permita recopilar noticias de distintas fuentes de información, y mediante el análisis automático de su contenido determine si este guarda relación con las secciones de interés del usuario y el periodo establecido. Finalmente, las noticias que satisfagan ambos filtros serán las que se le mostrarán al usuario.

## 1.2. Solución Propuesta

### 1.3. Objetivo

Crear un recolector de noticias, el cual permita recopilar información de diferentes fuentes como diarios, sitios de noticias, foros y mediante el análisis automático de su contenido muestre aquellas noticias que satisfagan los filtros de período y secciones establecidos por el usuario.

### 1.4. Objetivos Específicos

- Desarrollar un recolector de noticias, el cual permita recopilar información de diferentes fuentes como diarios, sitios de noticias, blogs y foros
- Analizar de forma automática el contenido de las noticias para satisfacer los filtros establecidos por el usuario
- Mostrar el enlace (URL) de las noticias que cumplieron con los filtros establecidos
- Afinar el clasificador de noticias realizado en el trabajo terminal 2017-A02 para utilizarlo en el contexto de esta propuesta (filtro de sección)

## **1.5. Estructura del Documento**

1.5. ESTRUCTURA DEL DOCUMENTO CAPÍTULO 1. INTRODUCCIÓN



## Capítulo 2

# Conclusiones

Vestibulum faucibus nibh ac felis dapibus, ac efficitur ipsum fermentum. Nullam sapien ligula, varius sed neque nec, dictum scelerisque ipsum. Praesent pellentesque tristique lorem non cursus. In et urna lectus. Cras porttitor ipsum sed ullamcorper faucibus. Curabitur sapien turpis, vulputate sed enim a, feugiat aliquet eros. Donec ut dui a libero dapibus dictum ac eget justo. Duis tristique luctus diam, faucibus eleifend neque tincidunt at. Sed eget risus dolor. Quisque auctor tellus eget ipsum maximus, at sodales nisi maximus. Fusce nisi lectus, ornare sit amet mi et, fermentum vestibulum turpis. Nulla et rhoncus nulla. Suspendisse arcu ligula, mollis sed diam a, consequat sollicitudin ipsum. Donec eget aliquet nisi.



## Capítulo 3

# Bibliografía

### 3.1. Referencias

(1) El Economista, Ranking de medios nativos digitales", Disponible en:  
<https://www.eleconomista.com.mx/Ranking-de-Medios-Nativos-Digitales>