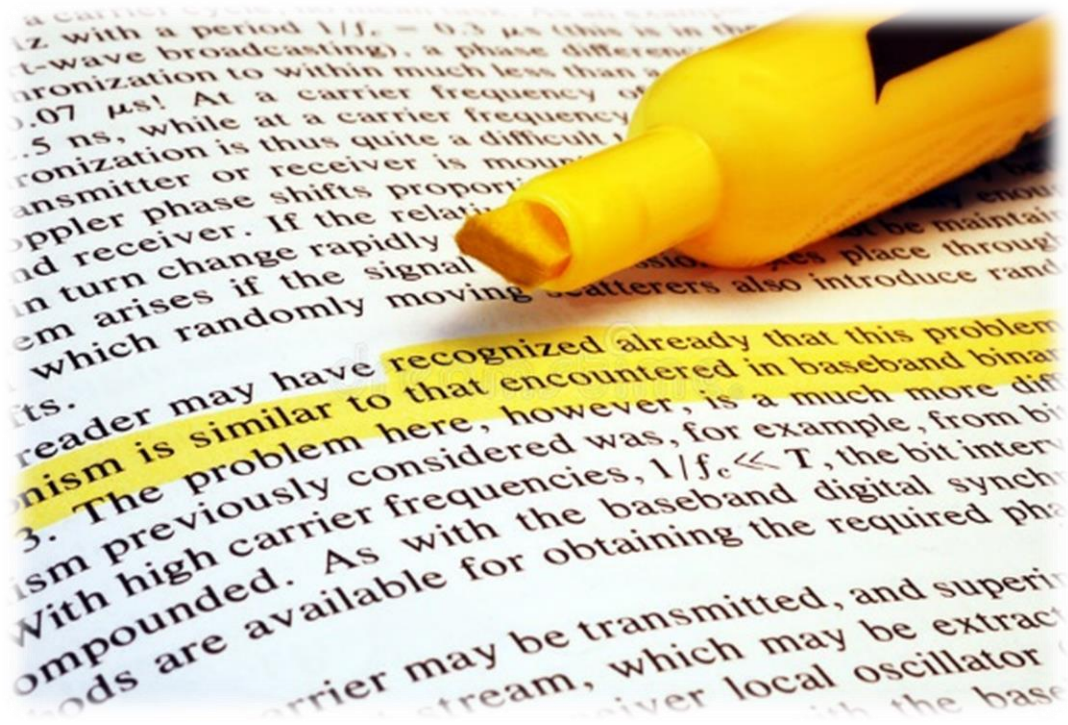




UTEM

UNIVERSIDAD
TECNOLÓGICA
METROPOLITANA

Proyecto Semestral: Analizador de textos con el algoritmo de Porter



Daniel Vásquez Robles.
Mayo 2019.

Universidad Tecnológica Metropolitana.
Depto. Ingeniería en Informática.
Teoría de Autómatas.
Profesor Ricardo Corbinaud.

Tabla de Contenidos

ii

| | |
|---|----|
| Capítulo 1 Introducción e información general | 1 |
| ¿Qué es el algoritmo de Porter? | 2 |
| Text-Mining | 3 |
| Capítulo 2 Proyecto semestral de teoría de autómatas..... | 5 |
| Autómatas | 5 |
| Expresiones regulares | 5 |
| Funciones relevantes | 6 |
| Comparación de resultados obtenidos por el algoritmo..... | 7 |
| Conclusiones | 8 |
| Lista de referencias | 9 |
| Apéndice | 10 |
| Pronombres: | 10 |
| Artículos:..... | 44 |
| Prefijos: | 48 |
| Sufijos: | 56 |

Lista de tablas

iii

| | |
|--|---|
| Tabla 1. Comparación de resultados obtenidos | 6 |
|--|---|

Capítulo 1

Introducción e información general

En el lenguaje cotidiano muchas veces utilizamos palabras que tratan del mismo concepto, sin embargo, estas cambian ligeramente dependiendo en el contexto en el cual sean aplicadas. Pues bien, en este proyecto se trabajó precisamente con esta idea basándonos en un algoritmo publicado bajo el nombre de “*An algorithm for suffix stripping*” por Martin F. Porter en el año 1980, Cambridge (Inglaterra).

En las siguientes paginas podrá encontrar una breve explicación de que es el algoritmo de Porter, como funciona y un poco la historia de este. Seguido, podrá encontrar información al respecto del *text-mining* concepto ligado a la recuperación de información dentro de la informática. Y para finalizar se mostrará el trabajo realizado a lo largo del semestre tratando los autómatas, expresiones regulares y funciones utilizadas en el desarrollo del algoritmo.

¿Qué es el algoritmo de Porter?

El algoritmo de derivación de Porter es un proceso para eliminar las terminaciones morfológicas e inflexionales más comunes de las palabras en inglés, sin embargo, para efectos de la realización de este proyecto se aplicará a las palabras en español.

El documento original del algoritmo de derivación se escribió en 1979 en el laboratorio de computación, Cambridge (Inglaterra), como parte de un proyecto de RI (recuperación de información) más grande, el cual sería publicado en 1980 con el nombre de “*An algorithm for suffix stripping*”.

Originalmente el algoritmo fue escrito en BCPL, un lenguaje que alguna vez fue popular, pero que ahora está extinguido. Durante los primeros años posteriores a su publicación, se distribuyó en su versión BCPL, a través de una cinta de papel perforada. Pronto comenzaron a aparecer versiones en otros lenguajes, y para 1999 se estaba utilizando, citando y adaptando ampliamente.

Remover los sufijos de las palabras automáticamente significa una operación que resulta especialmente útil en el campo de la recuperación de información. Podemos decir que un documento está formado por un vector largo de términos, los cuales pueden poseer una raíz en común, por ejemplo:

Conectar Conectado Conectando Conexión Conexiones

Para este ejemplo podemos dar cuenta que si a estas palabras se les extrae el sufijo la raíz en común de estas palabras es “conec-”, de esta forma el algoritmo realiza un conteo de

la frecuencia de las diferentes raíces obtenidas y con esta información podremos deducir de que trata el documento analizado por este algoritmo.

Algoritmos basados en este método conocido como Stemming (reducir una palabra a su raíz o stem) ayudan en sistemas de recuperación de información ya que aumentan el *recall* que es una medida sobre el número de documentos que se pueden encontrar con una consulta. Por ejemplo: Al obtener el stem “bibliotec-“, encontrara documentos para “biblioteca” y “bibliotecario” ya que ambas palabras poseen dicho stem en común.

El motor de búsqueda *Muscat* proviene de una investigación realizada por Porter en la Universidad de Cambridge y fue comercializado en 1984 por Cambridge CD Publishing; posteriormente se vendió a MAID, que se convirtió en la Corporación Dialog.

Text-Mining

Consiste en el proceso de analizar colecciones de materiales textuales con el fin de capturar conceptos y temas clave, descubrir relaciones y tendencias ocultas sin requerir que conozca las palabras o términos precisos que los autores han usado para expresar esos conceptos. Se suele confundir el text-mining con la recuperación de información, sin embargo, tienen un objetivo diferente ya que el text-mining se centra en la extracción y gestión de contenido de calidad, la terminología y las relaciones contenidas en la información.

En el text-mining destacan tres técnicas. La primera de ellas es la extracción de términos, la cual identifica los términos clave y entidades lógicas. Luego la extracción de información, consiste en identificar las relaciones básicas entre los términos extraídos del

texto. Por último, el análisis relacional permite formar modelos complejos, que permiten dar una idea de las relaciones entre varias entidades.

Al utilizar text-mining se pueden analizar grandes volúmenes de información de manera más rápida y eficiente en comparación a la codificación manual, ya que una vez establecidos y validados los parámetros, el proceso es automatizado. Además, facilita el análisis de datos no estructurados permitiendo así filtrar, buscar y hacer referencias cruzadas de manera más sencilla.

Capítulo 2

Proyecto semestral de teoría de autómatas

Autómatas

En el primer avance de proyecto se trabajó modelando autómatas que permitieran satisfacer los requerimientos exigidos. Como resultado se obtuvieron los diferentes autómatas para reconocer las diferentes categorías de palabras ya sean los sufijos, prefijos, artículos y pronombres. Cabe destacar que el modelo de autómatas es el mismo tratándose de la misma categoría de palabra, es decir, se puede aplicar el autómata tanto para un sufijo como para otro sufijo diferente, lo mismo aplica para cada categoría anteriormente mencionada. Todas las figuras correspondientes a los autómatas están ubicadas en el apéndice.

Expresiones regulares

Posteriormente se trabajó en el segundo avance de proyecto, esta vez modelando expresiones regulares capaces de identificar los sufijos, prefijos, pronombres o artículos presentes en un texto. En el caso de los sufijos se modeló por ejemplo que las palabras que contenían sufijos primero venía una serie de combinaciones de letras y luego le seguía el sufijo en cuestión. Para los prefijos, caso contrario, la palabra iniciaba con el prefijo y le seguía una combinación de letras.

Previo análisis del texto se utilizaron otras expresiones regulares para normalizar el documento, es decir, se transformó el texto solo a minúsculas, se eliminó los símbolos, tildes y números. De esta manera el algoritmo podrá cumplir su función.

Funciones relevantes

Finalmente, luego de trabajar con los avances, se sintetizó las ideas en un código JAVA, que como resultado de analizar un documento .txt crea un “Diccionario.txt” que contiene los 10 stem (o raíces) más frecuentes. Además, se observó que las expresiones regulares (regex) obtenidos en el avance dos, encapsulaban la palabra en su totalidad por lo que de hacer esto, eliminaría la palabra completa generando resultados erróneos, es por esto que fue necesario re-definir los regex a utilizar por el código final.

A continuación, destacare algunas de las funciones del código que ayudaron a que el proyecto cumpliera con su objetivo:

La función “Tildes” recibe un tipo de dato string que contiene una palabra obtenida del texto por el algoritmo y reemplaza los caracteres de vocales que presenten tildes por su formato sin tilde, esta es una de las funciones que ayuda a normalizar el texto como se mencionó anteriormente.

```
public String Tildes(String s){
    s = s.replaceAll("[àâä]", "a");
    s = s.replaceAll("[éèë]", "e");
    s = s.replaceAll("[îï]", "i");
    //s = s.replaceAll("ñ", "n"); //En caso de querer sustituir las "ñ" por "n".
    s = s.replaceAll("[óòö]", "o");
    s = s.replaceAll("[úùü]", "u");
    s = s.replaceAll("[ýÿ]", "y");
    return s;}
```

La siguiente función “Articulos” recibe un tipo de dato string que contiene una palabra obtenida del texto, verifica si es un artículo especificado por medio de la expresión regular, y en caso de que así sea, lo elimina. Esto lo realizamos por que los artículos se repiten mucho a lo largo de los textos y de no eliminarlos en los resultados aparecerían entro de los diez primeros stem más frecuentes entorpeciendo nuestro resultado final ya que los artículos no nos sugieren de que trata el texto.

```
public String Articulos(String s){
    s = s.replaceAll("\\b[en|el|l[o|a](s)?]\\b", "");
    return s;}

```

Comparación de resultados obtenidos por el algoritmo

A continuación, mostrare los resultados obtenidos por el algoritmo aplicando primero la eliminación de los sufijos y prefijos de las palabras, para luego comparar los resultados obtenidos solo aplicando la eliminación de sufijos.

Tabla 1. Comparación de resultados obtenidos.

| <i>Sufijos y Prefijos</i> | <i>Solo Sufijos</i> |
|--|--|
| 'utomata' se ha repetido 112 veces. | 'automata' se ha repetido 112 veces. |
| 'puede' se ha repetido 78 veces. | 'como' se ha repetido 112 veces. |
| 'datos' se ha repetido 75 veces. | 'puede' se ha repetido 78 veces. |
| 'publicaciones' se ha repetido 71 veces. | 'datos' se ha repetido 75 veces. |
| 'estados' se ha repetido 68 veces. | 'publicaciones' se ha repetido 71 veces. |
| 'cada' se ha repetido 64 veces. | 'estados' se ha repetido 68 veces. |
| 'utomatas' se ha repetido 64 veces. | 'automatas' se ha repetido 64 veces. |
| 'utores' se ha repetido 63 veces. | 'cada' se ha repetido 64 veces. |
| 'fuente:' se ha repetido 60 veces. | 'autores' se ha repetido 63 veces. |
| 'cendios' se ha repetido 55 veces. | 'fuente:' se ha repetido 60 veces. |
| 'partir' se ha repetido 45 veces. | 'incendios' se ha repetido 55 veces. |
| 'forestales' se ha repetido 41 veces. | 'partir' se ha repetido 45 veces. |
| 'matriz' se ha repetido 41 veces. | 'matriz' se ha repetido 41 veces. |
| 'medi' se ha repetido 40 veces. | 'forestales' se ha repetido 41 veces. |
| 'modelo' se ha repetido 39 veces. | 'medi' se ha repetido 40 veces. |
| 'mayor' se ha repetido 38 veces. | 'modelo' se ha repetido 39 veces. |
| 'cendio' se ha repetido 37 veces. | 'mayor' se ha repetido 38 veces. |
| 'ejemplo' se ha repetido 36 veces. | 'entre' se ha repetido 38 veces. |
| 'forma' se ha repetido 35 veces. | 'incendio' se ha repetido 37 veces. |
| 'njunto' se ha repetido 35 veces. | 'ejemplo' se ha repetido 36 veces. |

Conclusiones

Como conclusión luego de comparar los resultados obtenidos en la tabla, los stem resultantes presentan pequeñas diferencias debido a que se aplicó un leve cambio al “filtro de palabras” por así decirlo, sin embargo, al solamente eliminar los sufijos la información obtenida resulta más fácil de deducir, ya que, por ejemplo, el stem “Presley” resulta más fácil de comprender, nos dice más en definitiva que el stem “sley”, de modo que al observar el ligero cambio en los resultados da a pensar que el algoritmo funciona en ambos casos pero para facilitar la deducción y comprensión de los resultados, es más practico solamente eliminar los sufijos.

Lista de referencias

- Martin F. Porter (2006). Recuperado de <https://tartarus.org/martin/PorterStemmer/index.html>
- Porter, M. F. (1980). “*An algorithm for suffix stripping*”. Recuperado de <https://www.cs.odu.edu/~jbollen/IR04/readings/readings5.pdf>
- Universidad de Alcalá (2018). Recuperado de <https://www.master-data-scientist.com/en-que-consiste-el-text-mining/>
- Rochina, Paula (2017). *¿Qué es y cuáles son las aplicaciones del text-minig?*. Recuperado de <https://revistadigital.inesem.es/informatica-y-tics/text-mining/>

Apéndice

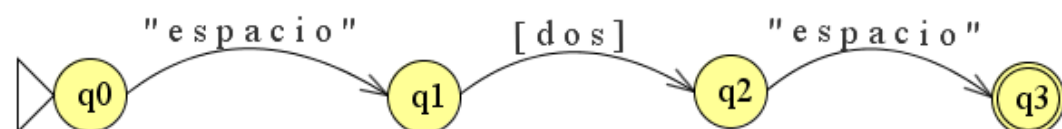
Pronombres:

El regex correspondiente es: $^{\wedge}[\text{pronombres}]\$$



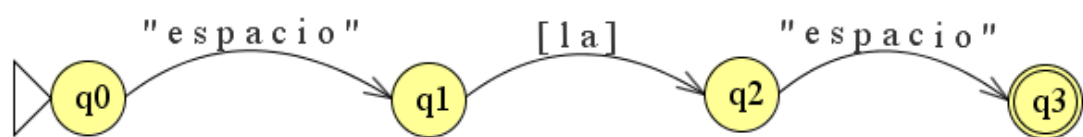




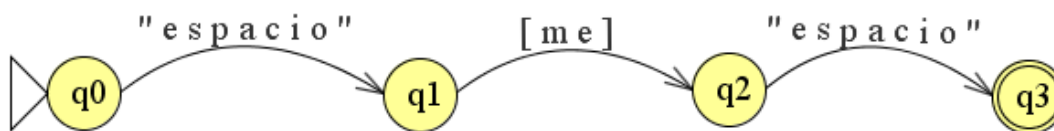
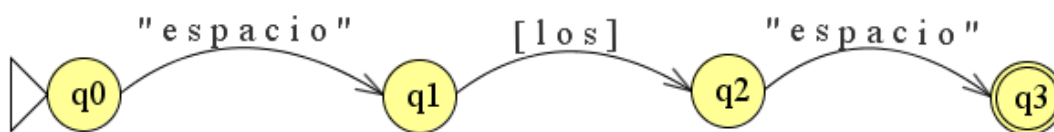
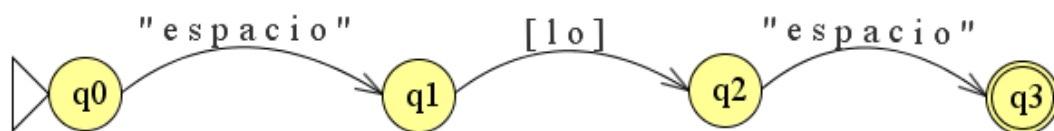


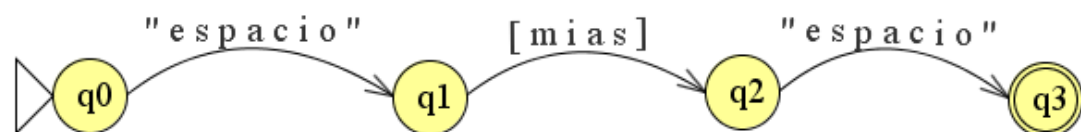
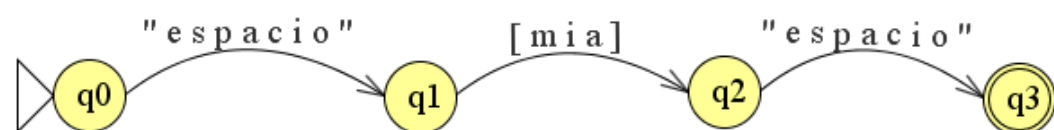
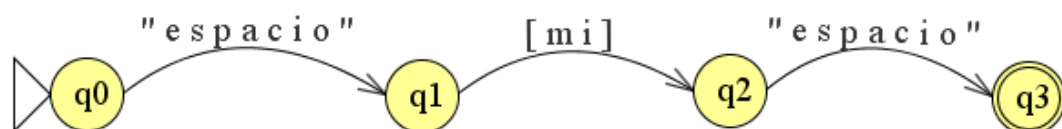


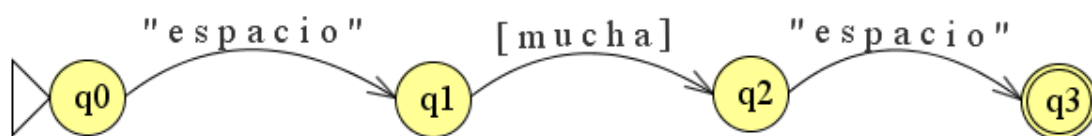
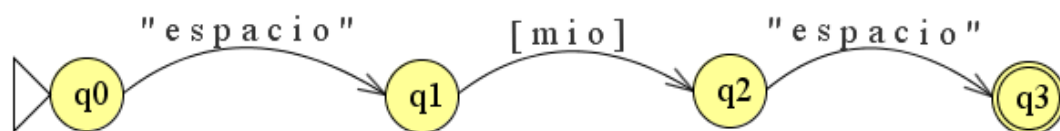


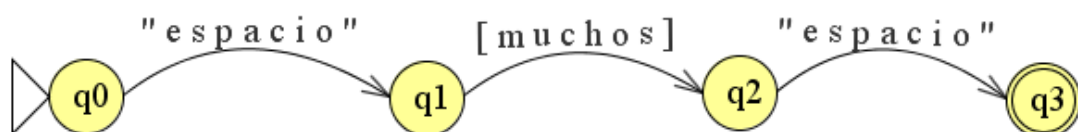
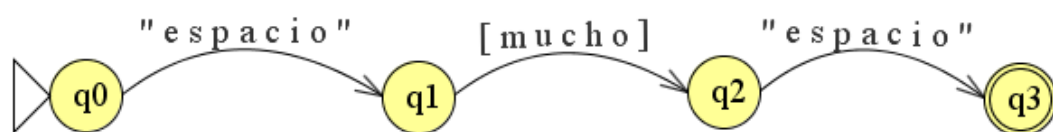
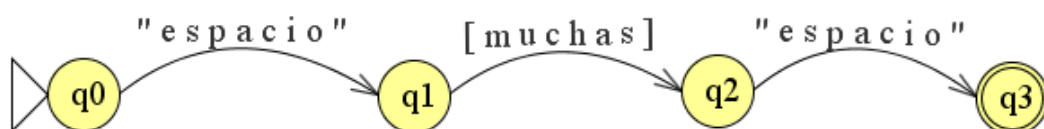


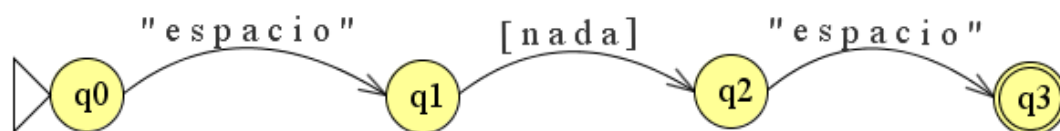


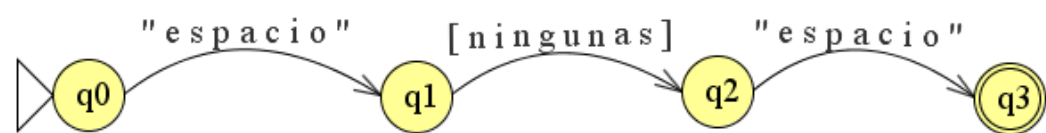








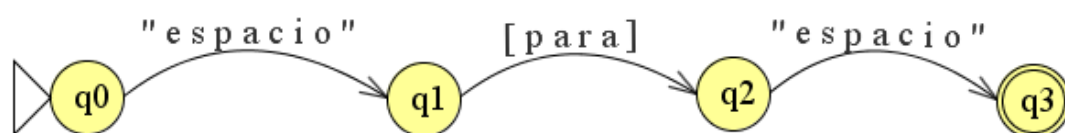






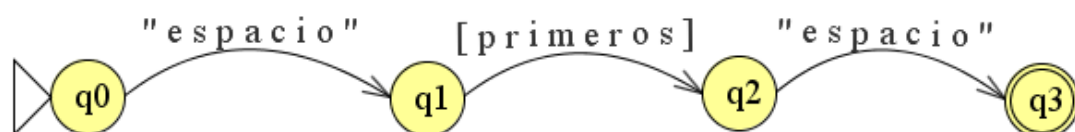




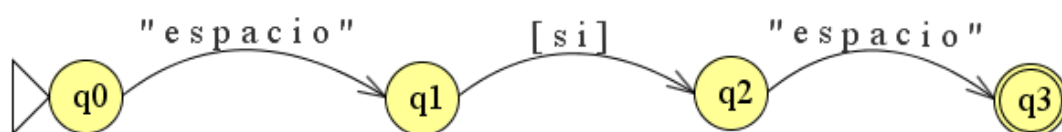


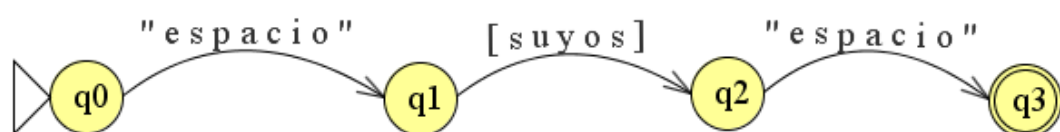
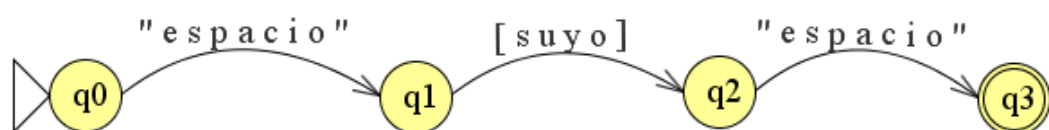
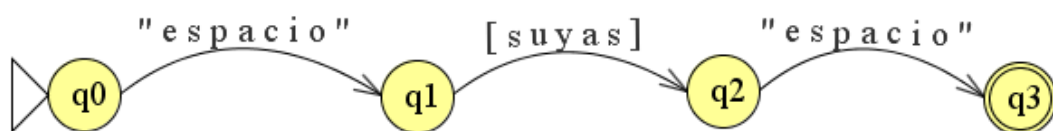








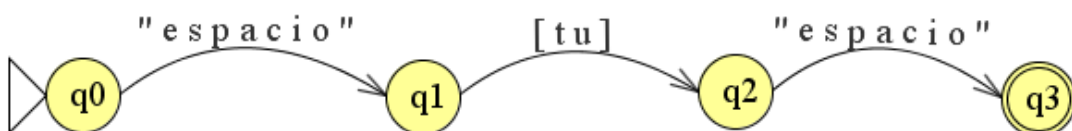


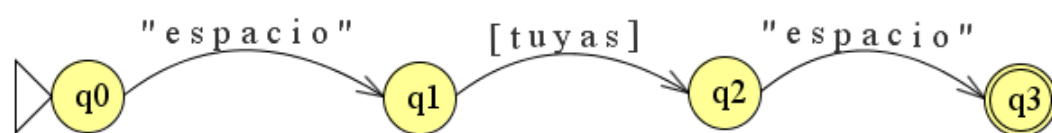
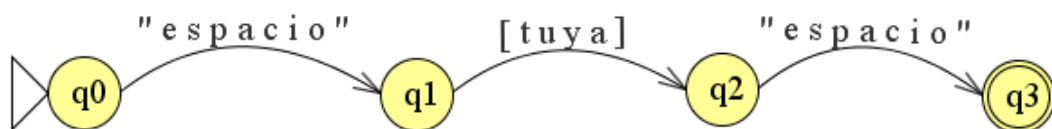


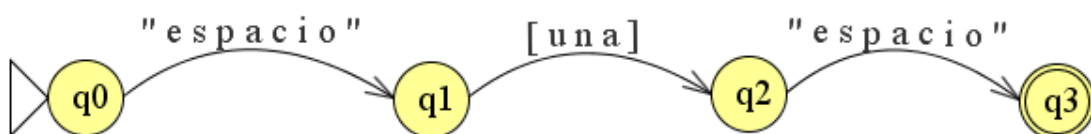
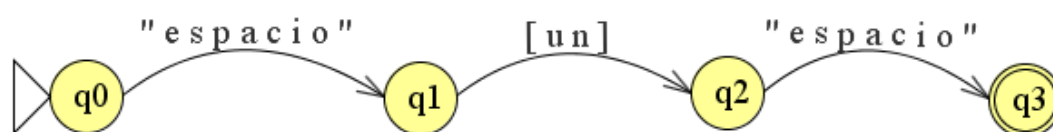
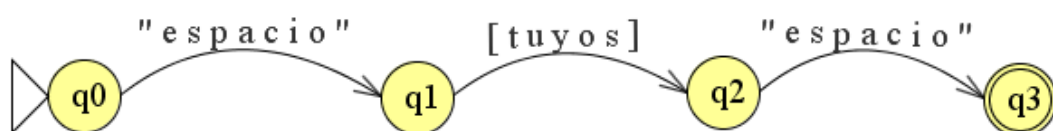


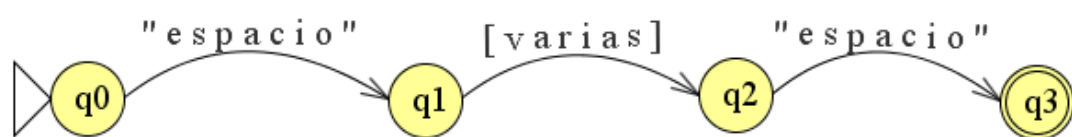
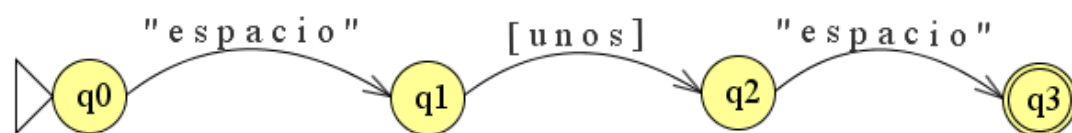
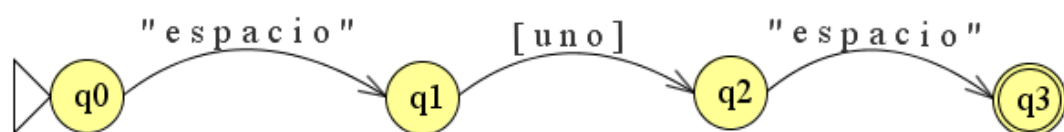


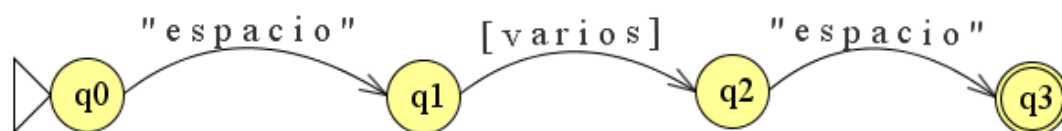




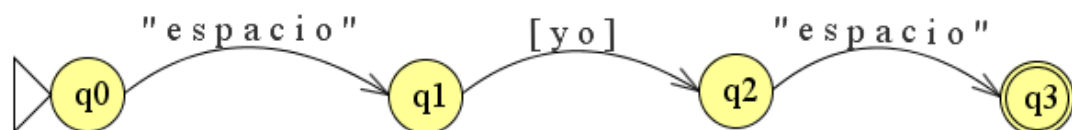






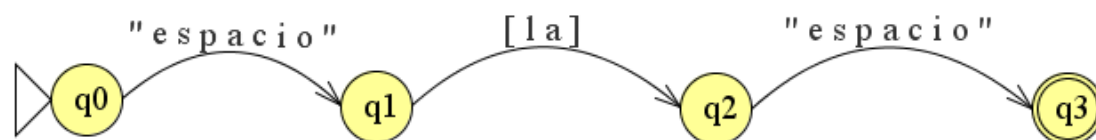
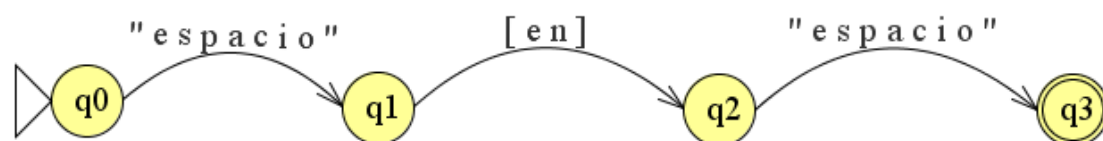
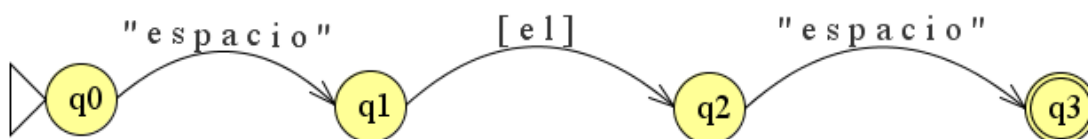


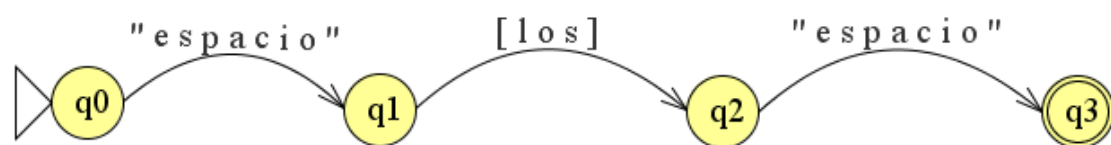
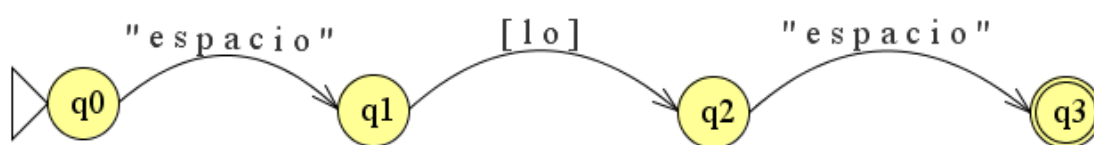
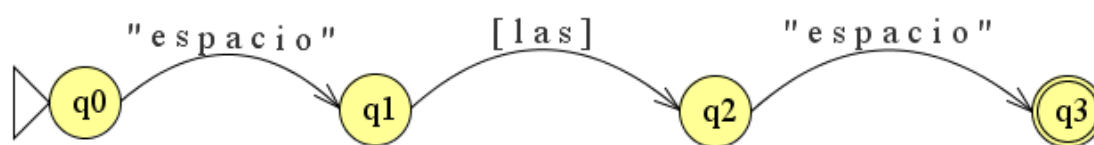


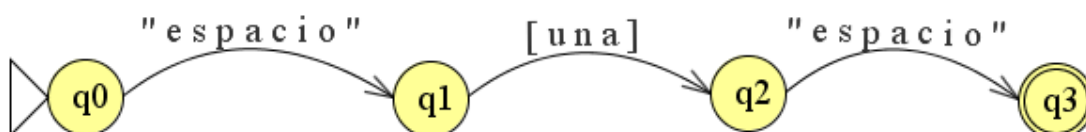
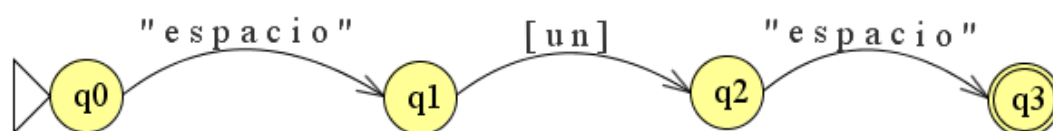


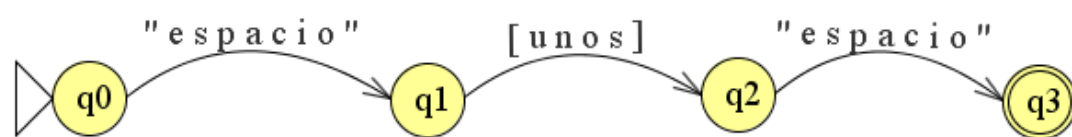
Articulos:

El regex correspondiente es: `^[articulo]$`



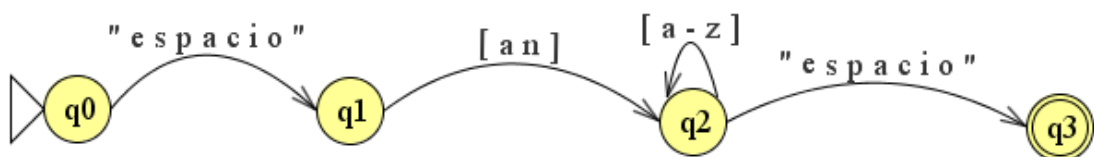
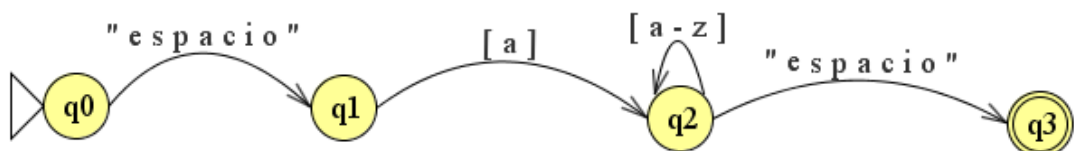


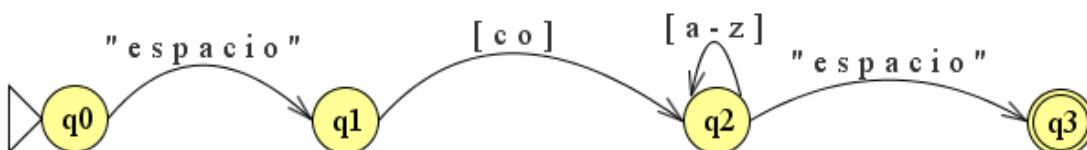
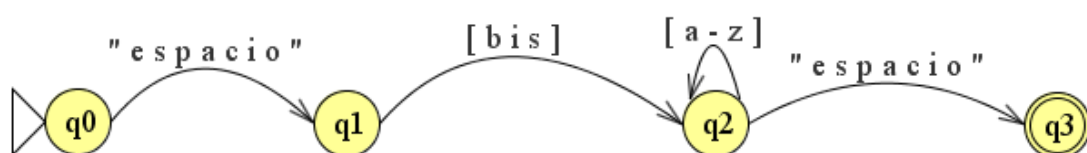
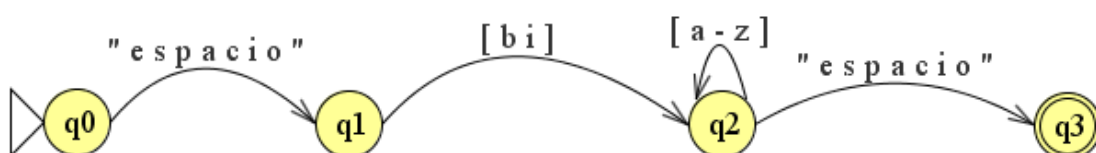
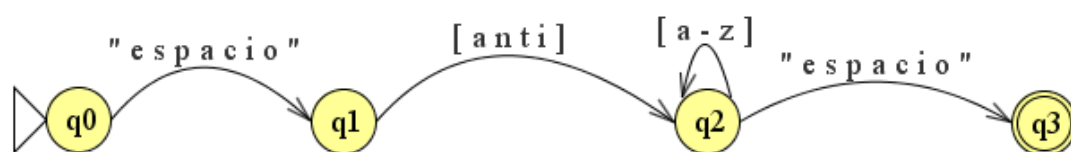


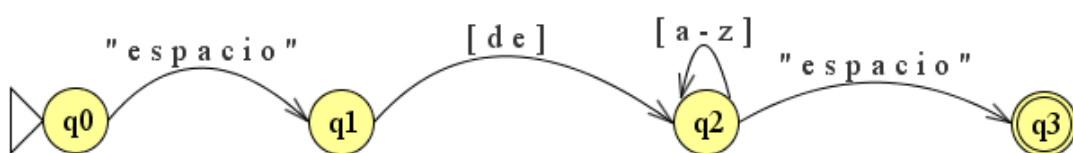


Prefijos:

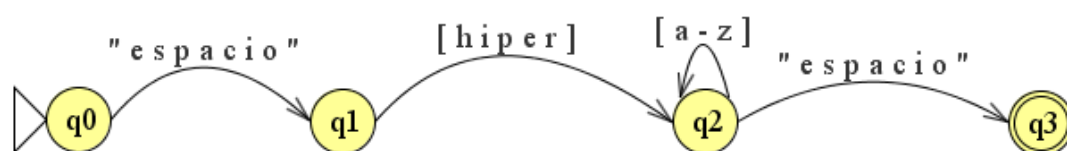
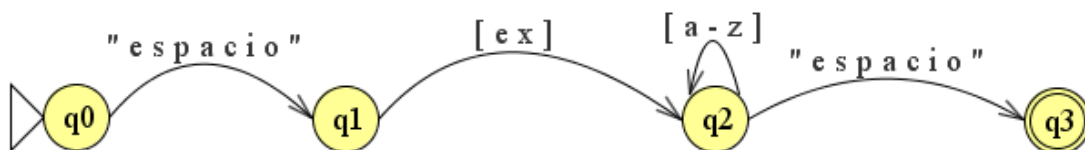
El regex correspondiente es: ^prefijo

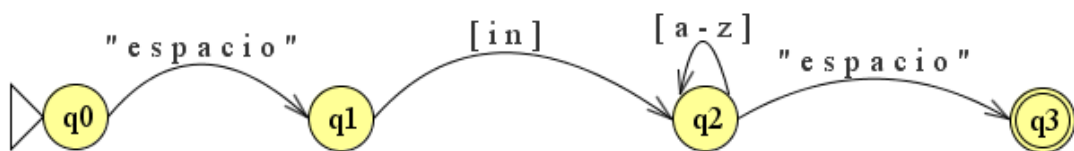
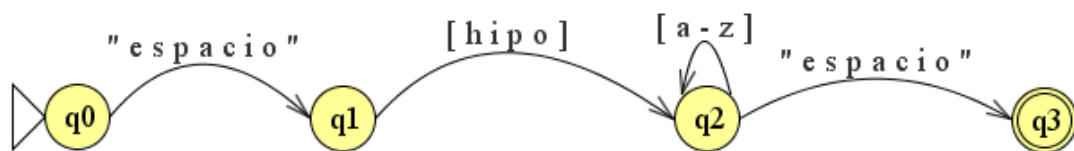


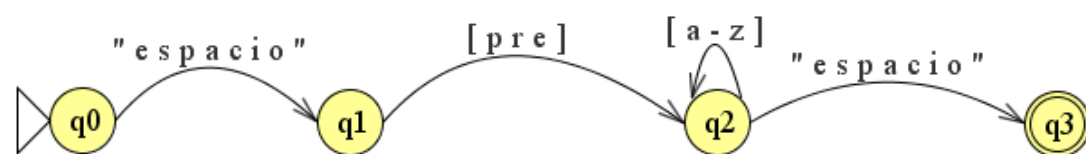
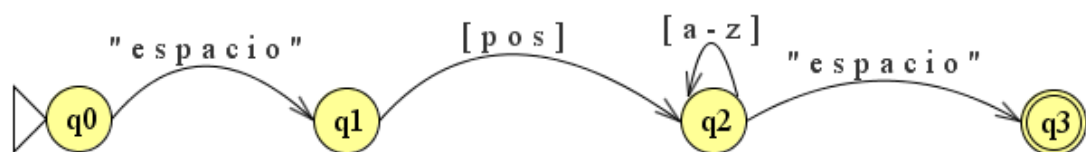


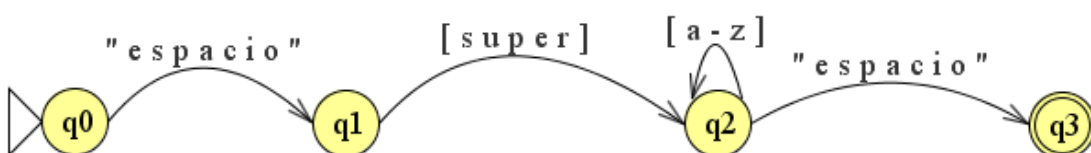
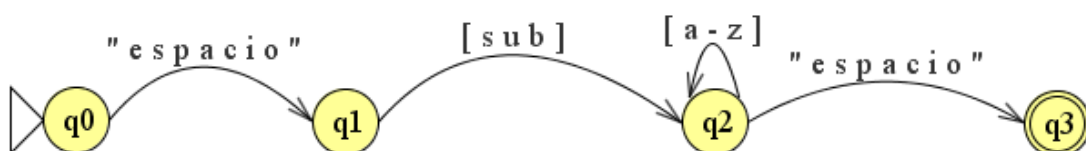
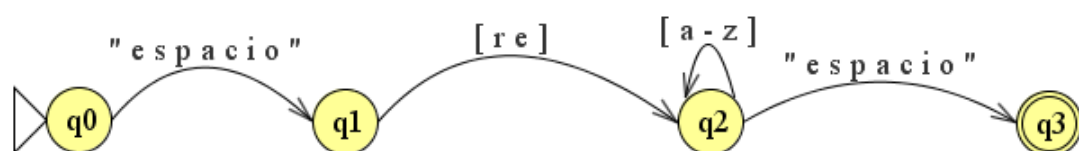












Sufijos:

El regex correspondiente es: sufijo\$

