

Expresiones regulares - grep

***Arquitectura y Sistemas Operativos
Tecnicatura Superior en Programación.
UTN-FRA***

Autores: *Prof. Martín Isusi Seff*

Revisores: *Prof. Marcos Pablo Russo*

Versión: 1



Esta obra está bajo una [Licencia Creative Commons Atribución-CompartirIgual 4.0 Internacional](https://creativecommons.org/licenses/by-sa/4.0/).

1. A continuación, utilizará el comando **wget** para descargar el contenido del sitio web de google.

```
alumno@alumno-VirtualBox:~/Google$ wget www.google.com
--2017-08-23 15:01:01-- http://www.google.com/
Resolving www.google.com (www.google.com)... 64.233.191.105, 64.233.191.103, 64.233.191.104, ...
Connecting to www.google.com (www.google.com)|64.233.191.105|:80... connected.
HTTP request sent, awaiting response... 302 Found
Location: http://www.google.com.ar/?gfe_rd=cr&ei=3sKdWdT-EIzgCLqNo9AB [following]
--2017-08-23 15:01:03-- http://www.google.com.ar/?gfe_rd=cr&ei=3sKdWdT-EIzgCLqNo9AB
Resolving www.google.com.ar (www.google.com.ar)... 64.233.191.94
Connecting to www.google.com.ar (www.google.com.ar)|64.233.191.94|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: unspecified [text/html]
Saving to: 'index.html'

index.html           [ <=>          ] 10,54K --.-KB/s   in 0,002s

2017-08-23 15:01:05 (5,77 MB/s) - 'index.html' saved [10797]

alumno@alumno-VirtualBox:~/Google$ ls
index.html
alumno@alumno-VirtualBox:~/Google$
```

Como se puede apreciar en la imagen, se utilizó el comando **wget www.google.com** y en la salida del mismo se ve que se realiza una conexión y se "salva" el archivo **index.html**. Podemos confirmar esto ejecutando el comando **ls**. Allí vemos que efectivamente el archivo está en el directorio.

2. Revisando el contenido del archivo, vemos que el mismo no es más ni menos que código HTML¹. En HTML, cada elemento de una página web (imágenes, links, párrafos, etc.) se representa con etiquetas de la siguiente manera:

<etiqueta>

Contenido del elemento.

</etiqueta>

Nótese que el contenido del elemento está "encerrado" por etiquetas que se abren y se cierran. Existen distintos tipos de etiquetas. Por ejemplo, los links se representan de la siguiente manera:

<a>Título

Dentro de las etiquetas se puede agregar más información. Esto da, por ejemplo, que tengamos etiquetas de la siguiente manera:

Link

Analizando la estructura de HTML, podemos decir que los links empiezan con **<a** y terminan con **/a>**, pudiendo haber cualquier carácter en el medio.

3. Habiendo completado el punto 1 y entendido el punto 2, en este punto se pide desarrollar una expresión regular que sirva para detectar links dentro del archivo descargado previamente. Una vez encontrados los links, redirigirlos hacia un nuevo archivo.

¹ Lenguaje de marcado utilizado para el maquetado de sitios web.