

Capítulo 10

Expresiones regulares

?	<i>El elemento precedente es opcional y debe coincidir al menos una vez.</i>
*	<i>El elemento precedente debe coincidir cero o más veces.</i>
{n}	<i>El elemento precedente debe coincidir exactamente n veces.</i>
+	<i>El elemento precedente debe coincidir una o más veces.</i>
{,m}	<i>El elemento precedente es opcional y debe coincidir al menos m veces.</i>
{n,m}	<i>El elemento precedente debe coincidir al menos n veces pero no más de m veces.</i>
pablo	<i>La cadena pablo.</i>
^pablo	<i>La cadena pablo al principio de la línea.</i>
pablo\$	<i>La cadena pablo al final de la línea.</i>
^pablo\$	<i>La cadena pablo formando una única línea.</i>
niñ[oa]	<i>La cadena niño o niña.</i>
ni[^aeiou]	<i>La tercera letra no es una vocal minúscula.</i>
ga.o	<i>La tercera letra es cualquier carácter.</i>
^....\$	<i>Cualquier línea que contenga 4 caracteres.</i>
^.	<i>Cualquier línea que comienza por punto.</i>
^[^.]	<i>Cualquier línea que no comienza por punto.</i>
niños*	<i>niño, niños, niñoss, niñosss, etc</i>
“niño”	<i>niño entre comillas dobles.</i>
“*niño”*	<i>niño con o sin comillas dobles.</i>
[a-z][a-z]*	<i>una o más letras minúsculas.</i>
[a-z]+	<i>una o más letras minúsculas (sólo válido en algunas aplicaciones).</i>
[^0-9A-Z]	<i>cualquier carácter que no sea ni número ni letra mayúscula.</i>
[a-zA-Z]	<i>cualquier letra sea mayúscula o minúscula.</i>
[Ax5]	<i>cualquier carácter que sea A, x o 5.</i>
niño niña nada	<i>una de las tres palabras.</i>
(s arb)usto	<i>la palabra susto o arbusto.</i>
ga?t[oa]	<i>gato, gata, gasto, gaita, etc.</i>
\<ga	<i>cualquier palabra que empiece por ga.</i>
ño\>	<i>cualquier palabra que termine por ño</i>
\<niño\>	<i>la palabra niño</i>
o\{2,\}	<i>dos o más oes en una misma fila.</i>
\t	<i>Representa un tabulador.</i>
\r	<i>Representa el "retorno de carro" o "regreso al inicio" o sea el lugar en que la línea vuelve a iniciar.</i>
\n	<i>Representa la "nueva línea" el carácter por medio del cual una línea da inicio.</i>
\a	<i>Representa una "campana" o "beep" que se produce al imprimir este carácter.</i>
\e	<i>Representa la tecla "Esc" o "Escape"</i>
\f	<i>Representa un salto de página.</i>
\v	<i>Representa un tabulador vertical.</i>
\x	<i>Se utiliza para representar caracteres ASCII o ANSI si conoce su código. De esta forma, si se busca el símbolo de derechos de autor</i>

	<i>y la fuente en la que se busca utiliza el conjunto de caracteres Latin-1 es posible encontrarlo utilizando "\xA9".</i>
\u	<i>Se utiliza para representar caracteres Unicode si se conoce su código. "\u00A2" representa el símbolo de centavos. No todos los motores de Expresiones Regulares soportan Unicode. El .Net Framework lo hace, pero el EditPad Pro no, por ejemplo.</i>
\d	<i>Representa un dígito del 0 al 9.</i>
\w	<i>Representa cualquier carácter alfanumérico.</i>
\s	<i>Representa un espacio en blanco.</i>
\D	<i>Representa cualquier carácter que no sea un dígito del 0 al 9.</i>
\W	<i>Representa cualquier carácter no alfanumérico.</i>
\S	<i>Representa cualquier carácter que no sea un espacio en blanco.</i>
\A	<i>Representa el inicio de la cadena. No un carácter sino una posición.</i>
\Z	<i>Representa el final de la cadena. No un carácter sino una posición.</i>
\b	<i>Marca el inicio y el final de una palabra.</i>
\B	<i>Marca la posición entre dos caracteres alfanuméricos o dos no-alfanuméricos.</i>

Buscar una expresión regular

El comando **grep** toma una expresión regular de la línea de comandos, lee la entrada estándar o una lista de archivos, e imprime las líneas que contengan coincidencias para la expresión regular.

Opciones :

-c	<i>Modificar la salida normal del programa, en lugar de imprimir por salida estándar las líneas coincidentes, imprime la cantidad de líneas que coincidieron en cada archivo.</i>
-e PATRÓN	<i>Usar PATRÓN como el patrón de búsqueda, muy útil para proteger aquellos patrones de búsqueda que comienzan con el signo «-».</i>
-f ARCHIVO	<i>Obtener los patrones del archivo ARCHIVO</i>
-H	<i>Imprimir el nombre del archivo con cada coincidencia.</i>
-r	<i>Buscar recursivamente dentro de todos los subdirectorios del directorio actual.</i>
-i	<i>Buscar mayúsculas o minúsculas.</i>
-w	<i>Buscar palabra completa.</i>
-v	<i>Buscar lo contrario.</i>
-l	<i>Mostrar solo los archivos que contenga la palabra buscada.</i>
-c	<i>Contar la cantidad de ocurrencia.</i>
-E	<i>Expresiones regulares extendidas en vez de usar grep podemos usar egrep que es lo mismo, la diferencia es que egrep son expresiones regulares extendidas y no tenemos que pasar la opción -E como en grep.</i>
-o	<i>Muestra solamente lo que queremos buscar en la línea.</i>

Ejemplo :

```
# grep 'root' /etc/passwd
# grep -i 'root' /etc/passwd
# grep -iw 'root' /etc/passwd
```

```
# grep -v 'root' /etc/passwd
# grep -l 'root' /etc/* | more
# grep -c 'root' /etc/* | more
```

Si queremos que tenga color la búsqueda realizamos la siguiente exportación :

```
# export GREP_OPTIONS='--color=auto' GREP_COLOR='100;8'
```

Modificar archivos, buscar, etc

El comando **sed** nos permite, de una forma cómoda, borrar líneas, registros o sustituir cadenas de caracteres dentro de las líneas, etc.

Opciones :

-n, --quiet, --silent	<i>Suprime la muestra automática del espacio de patrones.</i>
-e guión, --expresion=guión	<i>Agrega el guión a la lista de órdenes para ejecutar.</i>
-f fichero-guión, --file=fichero-guión	<i>Agrega el contenido del fichero guión a la lista de órdenes para ejecutar.</i>
--follow-symlinks	<i>Selecciona solamente una columna para mostrar.</i>
-i[SUFIJO], --in-place[=SUFIJO]	<i>Edita ficheros en el lugar (crea un respaldo si se da una extensión).</i>
-l N, --line-length=N	<i>Especifica la longitud de corte de línea deseado para la orden `l`.</i>
--posix	<i>Desactiva todas las extensiones de GNU.</i>
-r, --regexp-extended	<i>Utiliza expresiones regulares extendidas en el guión.</i>
-s, --separate	<i>Considera los ficheros como separados en lugar de un solo flujo, largo y continuo.</i>
-u, --unbuffered	<i>Carga cantidades mínimas de datos de los ficheros de entrada y vacía los almacenamientos temporales de salida con más frecuencia.</i>
--help	<i>Ayuda.</i>
--version	<i>Versión.</i>

Comandos :

d	<i>Borrar linea.</i>
i	<i>Insertar.</i>
s	<i>Sustituir.</i>
p	<i>Listar.</i>
\+	<i>Mas de una expresión.</i>
N	<i>Fuerza la lectura de la siguiente línea en la iteración actual.</i>
y	<i>Sustituir.</i>

Ejemplo :

```
# cp /etc/passwd .
# sed '1d' passwd | more
# sed '1,5d' passwd | more
```

```
# sed '1i Hola Mundo' passwd | more
# sed 's/root/pepe/' passwd | more
# sed -n '1,4p' passwd
# sed -n '/root/p' passwd
```

```
o
```

```
# grep 'root' passwd
```

Eliminar líneas en blanco

```
# sed '/^$/d' passwd
```

Ejemplo de obtener los mail de un archivo :

Explicaremos la expresión regular : `\b[A-Z0-9._%~+@-]+\.[A-Z]{2,4}` :

<code>\b[A-Z0-9._%~+@-]+</code>	Buscamos un comienzo de palabra <code>\b</code> que contenga uno o más caracteres + dentro del conjunto A a Z, 0 a 9 , <code>.</code> , <code>_</code> , <code>%</code> y <code>-</code> , que son los caracteres admitidos para un nombre de usuario.
<code>@</code>	Seguido del símbolo <code>@</code> .
<code>[A-Z0-9.-]+</code>	Seguido por uno o más + caracteres A a Z, 0 a 9 , <code>.</code> y <code>-</code> , que son los caracteres admitidos para un dominio.
<code>\.</code>	Seguido por un punto escapado con <code>\</code> , ya que el punto solo significa cualquier carácter.
<code>[A-Z]{2,4}</code>	Seguido por entre 2 y 4 caracteres <code>{2,4}</code> entre el conjunto A-Z y además éste tiene que ser el final de una palabra <code>\b</code> .

```
# echo "El mail de juan es juan@hotmail.com es programador" > mail.txt
# echo "El mail de pedro es pedro.al@hotmail.com ar es webmaster" >> mail.txt
```

```
# egrep -oi '\b[A-Z0-9._%~+@-]+\.[A-Z]{2,4}' mail.txt
```

```
juan@hotmail.com
pedro.al@hotmail.com
```

Obtener las ip de un archivo

```
# echo "ip del servidor 192.168.2.10" > mis-ip.txt
# echo "ip del router 10.1.1.1" >> mis-ip.txt
# echo "ip ana 192.168.2.100" >> mis-ip.txt
# echo "ip miguel 192.168.2.101" >> mis-ip.txt
# echo "ip pedro 192.168.2.102" >> mis-ip.txt
```

```
# egrep -oi '\b[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3}\b' mis-ip.txt
```

```
192.168.2.10
10.1.1.1
```

```
192.168.2.100
192.168.2.101
192.168.2.102
```

```
o
```

```
# grep -oE '\b[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3}\b' mis-ip.txt
```

Listar los usuarios con password a partir del `/etc/shadow`

Resulta que queremos saber los usuarios del fichero `/etc/shadow` que pueden logearse al sistema (que no tienen un carácter `!` o `*` en el campo password). La expresión regular a usar sería:

<code>^[^:]+</code>	Buscamos uno o más caracteres + que no sean <code>:</code> a principio de línea <code>^</code>
<code>:</code>	Seguido de un carácter <code>:</code>
<code>[^*!]+</code>	Seguido por una secuencia de uno o más caracteres que no sean <code>!</code> o <code>*</code>
<code>:</code>	Seguido de otro carácter <code>:</code>

```
# egrep '^[^:]+:[^*!]+:' /etc/shadow
```

```
fulano:$1$PFC1F2pj$0KEfLTjNmDQg9.rtFT/DK0:12793:0:99999:7:::
mengano:$1$aAUa4js/$c2LEh.HeT0JvCG3B.340.:12793:0:99999:7:::
```

Borra desde `root` hasta `nobody` inclusive.

```
# sed '/root/,/nobody/d' passwd
```

Si queremos obtener un listado sólo de los directorios (las filas que comienzan con la letra "d"). Podríamos utilizar el sed para borrar las entradas que no comienzan con esa letra:

```
# ls -l | sed '/^d/!d'
```

Como vemos el símbolo de exclamación "`!`" detrás de la expresión regular niega la condición. Debería leerse como borrar todas las líneas que no comiencen con la letra '`d`'

Otra forma :

```
# ls -l | sed '/^d/p'
```

Sustituir los valores 123 por 456:

```
# echo "123" | sed 'y/[123]/[456]/'
```

```
456
```

Sacar los tabuladores de un archivo y reemplazarlo por un espacio :

```
# vi texto.txt
Hola Mundo ! !!!!
```

```
# cat texto.txt | sed 's/\t\+/ /g'
```

Sacar los mas de un espacio de un archivo y reemplazarlo por un espacio :

```
# vi texto.txt
  Hola Mundo !   !!!!

# cat texto.txt | sed 's/\+/ /g'
```

Procesando varias lineas

```
# vi agenda.txt

  Registro: 1
  Nombre: Anibal
  Telefono: 621-229

  Registro: 2
  Nombre: Hector
  Telefono: 562-245

  Registro: 3
  Nombre: Pablo
  Telefono: 622-354
```

Y queremos obtener lo siguiente :

```
1;Anibal;621-229
2;Hector;562-245
3;Pablo;622-354
```

Una solución es poder leer de a una línea a la vez. Este problema se resuelve con el comando “**N**”, que fuerza la lectura de la siguiente línea en la iteración actual.

```
# sed 'N;N;N;s/\n/;/g;s/^Registro: \(.*\);Nombre: \(.*\);Telefono: \(.*\);$/\1;\2;\3/'
agenda.txt
```

Son 3 líneas (N).

Que reemplace el retorno de carro “**\n**” por un “**;**”

```
buffer = "Registro: 1;Nombre: Anibal;Telefono: 621-229;"
```

Sacar las palabras “**Registro:** ”, “**Nombre:** ”, “**Telefono:** ” y el punto y coma de la linea.

Los paréntesis **\(.*)** representa a cada valor **\1**, **\2** y **\3**.

