

Práctica 2

Análisis de ficheros log de acceso a páginas web

12 de mayo de 2015

Toda la actividad que realizan los usuarios en un servidor web queda reflejada en sus ficheros *log*. Estos ficheros guardan información textual como la que aparece en la figura siguiente:

```
File Edit Options Buffers Tools Help
66.249.78.47 - - [30/Nov/2014:08:48:43 +0100] "GET /homepage/docs/CV_52182754D.pdf HTTP/1.1" 304 189 "-" "Mozilla/5.0 (compatible; Googlebot/2.1; +http://www.google.com/bot.html)"
217.176.218.87.dynamic.jazztel.es - - [30/Nov/2014:09:06:13 +0100] "GET /moodle/pluginfile.php/4578/user/icon/elegance/f2?rev=6114 HTTP/1.1" 200 3109 "-" "Mozilla/5.0 (Windows NT 6.3; Win64; x64; Trident/7.0; Touch; MSAppHost/2.0; rv:11.0) like Gecko"
217.176.218.87.dynamic.jazztel.es - - [30/Nov/2014:09:07:05 +0100] "GET /moodle/pluginfile.php/4637/user/icon/elegance/f2?rev=5965 HTTP/1.1" 200 3264 "-" "Mozilla/5.0 (Windows NT 6.3; Win64; x64; Trident/7.0; Touch; MSAppHost/2.0; rv:11.0) like Gecko"
crawl-66-249-64-61.googlebot.com - - [30/Nov/2014:09:11:58 +0100] "GET /moodle/theme/switchdevice.php?url=http%3A%2F%2Fwild.mat.ucm.es%2Fmoodle%2Flogin%2Fforgot_password.php&device=default&sesskey=A5bN6RgLXP HTTP/1.1" 404 5227 "-" "Mozilla/5.0 (iPhone; CPU iPhone OS 6_0 like Mac OS X) AppleWebKit/536.26 (KHTML, like Gecko) Version/6.0 Mobile/10A5376e Safari/8536.25 (compatible; Googlebot/2.1; +http://www.google.com/bot.html)"
crawl-66-249-78-33.googlebot.com - - [30/Nov/2014:09:40:40 +0100] "GET /moodle/calendar/set.php?return=L2NhbGVuZGFyL3ZpZxcucGhpP3ZpZxc9bW9udGmdGltZT0xNDUxNjAyODAwJmNvdXJzZT0x&sesskey=A68SuisWDW&var=showcourses HTTP/1.1" 404 5135 "-" "Mozilla/5.0 (compatible; Googlebot/2.1; +http://www.google.com/bot.html)"
crawl-66-249-78-40.googlebot.com - - [30/Nov/2014:09:47:39 +0100] "GET /moodle/theme/switchdevice.php?url=http%3A%2F%2Fwild.mat.ucm.es%2Fmoodle%2Fdevice=default&sesskey=Ns2CmBzR2 HTTP/1.1" 404 5226 "-" "Mozilla/5.0 (iPhone; CPU iPhone OS 6_0 like Mac OS X) AppleWebKit/536.26 (KHTML, like Gecko) Version/6.0 Mobile/10A5376e Safari/8536.25 (compatible; Googlebot/2.1; +http://www.google.com/bot.html)"
crawl-66-249-78-40.googlebot.com - - [30/Nov/2014:09:48:41 +0100] "GET /moodle/?lang=eu HTTP/1.1" 200 9346 "-" "Mozilla/5.0 (compatible; Googlebot/2.1; +http://www.google.com/bot.html)"
crawl-66-249-78-40.googlebot.com - - [30/Nov/2014:09:51:48 +0100] "GET /moodle/calendar/view.php?view=month&time=1417337321&course=1 HTTP/1.1" 200 7188 "-" "Mozilla/5.0 (compatible; Googlebot/2.1; +http://www.google.com/bot.html)"
google-proxy-66-249-93-155.google.com - - [30/Nov/2014:09:56:02 +0100] "GET /moodle/pluginfile.php/4578/user/icon/elegance/f2?rev=6114 HTTP/1.1" 200 3109 "-" "Mozilla/5.0 (Windows NT 5.1; rv:11.0) Gecko Firefox/11.0 (via ggpht.com GoogleImageProxy)"
crawl-66-249-78-33.googlebot.com - - [30/Nov/2014:09:56:17 +0100] "GET /moodle/calendar/view.php?view=month&time=1417337321&lang=ca HTTP/1.1" 200 7139 "-" "Mozilla/5.0 (compatible; Googlebot/2.1; +http://www.google.com/bot.html)"
google-proxy-66-249-93-155.google.com - - [30/Nov/2014:09:56:19 +0100] "GET /moodle/login/index.php HTTP/1.1" 303 909 "-" "Mozilla/5.0 (iPad; CPU OS 8_1_1 like Mac OS X) AppleWebKit/600.1.4 (KHTML, like Gecko) CriOS/39.0.2171.50 Mobile/12B436 Safari/600.1.4"
google-proxy-66-249-93-149.google.com - - [30/Nov/2014:09:56:19 +0100] "GET /moodle/login/index.php HTTP/1.1" 200 6271 "-" "Mozilla/5.0 (iPad; CPU OS 8_1_1 like Mac OS X) AppleWebKit/600.1.4 (KHTML, like Gecko) CriOS/39.0.2171.50 Mobile/12B436 Safari/600.1.4"
google-proxy-66-249-93-155.google.com - - [30/Nov/2014:09:56:19 +0100] "GET /moodle/theme/yui_combo.php?rollup/3.13.0/yui-moodlesimple-min.js HTTP/1.1" 304 307 "http://wild.mat.ucm.es/moodle/login/index.php" "Mozilla/5.0 (iPad; CPU OS 8_1_1 like Mac OS X) AppleWebKit/600.1.4 (KHTML, like Gecko) CriOS/39.0.2171.50 Mobile/12B436 Safari/600.1.4"
google-proxy-66-249-93-149.google.com - - [30/Nov/2014:09:56:19 +0100] "GET /moodle/theme/jquery.php/core/jquery-1.10.2.min.js HTTP/1.1" 304 10000 "http://wild.mat.ucm.es/moodle/theme/jquery.php/core/jquery-1.10.2.min.js" "Mozilla/5.0 (iPad; CPU OS 8_1_1 like Mac OS X) AppleWebKit/600.1.4 (KHTML, like Gecko) CriOS/39.0.2171.50 Mobile/12B436 Safari/600.1.4"
---- access.log.23 1% L?? (Fundamental)
Undo!
```

Figura 1: Fragmento de un fichero log.

Mientras que un libro suele tener entre 4000 y 6000 líneas, en un solo día, un servidor apache de un sitio 'modesto' como es `wild.mat.ucm.es`, puede generar decenas de miles de líneas. De forma simplificada, podemos decir que por cada petición que hace un usuario del servidor, se registra una línea que contiene toda la información relacionada con la petición.

Los ficheros log tienen un formato muy concreto y son muy fácilmente legibles (aunque mirando la imagen anterior cueste creerlo!). Miremos una única línea:

```
147.96.18.203 - - [27/Apr/2015:10:15:32 +0200] "GET /moodle/ HTTP/1.1" 200 8944 Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:34.0) Gecko/20100101 Firefox/34.0"
```

La información se estructura por columnas separadas por espacios en blanco. Vamos a comentar algunas columnas que serán de nuestro interés, para una información más

detallada puede consultarse la documentación oficial (<http://httpd.apache.org/docs/1.3/logs.html#accesslog>).

- En la primera columna aparece la dirección ip del cliente, lo que llamaremos usuarios. En el ejemplo anterior 147.96.18.203.
- En la cuarta columna aparece la fecha y la hora en la que la petición de página se ha realizado, [27/Apr/2015:10:15:32 +0200].
- En la quinta columna aparece la petición de página concreta que el usuario ha hecho, "GET /moodle/ HTTP/1.1". De esta petición lo que más nos interesa es la página pedida, en este caso /moodle/.
- En la sexta columna aparece un código de respuesta, que indica si se ha podido responder a la petición del usuario, si la página no existe, etc. (véase, por ejemplo, http://es.wikipedia.org/wiki/Anexo:C%C3%B3digos_de_estado_HTTP). En el ejemplo anterior el código de respuesta es 200, que indica que la página solicitada ha sido enviada sin ningún problema.

El análisis de log es una actividad que interesa a muy diversos niveles: identificar las páginas más visitadas, la carga del servidor en las distintas horas del día, los navegadores utilizados por los usuarios... Vamos a hacerte algunas propuestas concretas:

1 Sesiones

Decimos que una *sesión de usuario* es la actividad que ha realizado de forma continuada. Diremos que dos peticiones de un usuario están en diferentes sesiones si difieren en más de un determinado tiempo T . Este tiempo suele estar entre media hora o una hora, pero queremos que sea un parámetro de nuestra solución para poder cambiarlo en cualquier momento.

Escribe un programa map-reduce que, a partir de unos ficheros log y un tiempo de cierre de sesión T , devuelva el número total de sesiones de cada usuario. Utiliza librerías externas que te permitan manejar con facilidad las fechas que hay en el log.

2 Comportamiento de los usuarios

Supongamos que estamos interesados en analizar el *comportamiento* de los usuarios en cada sesión. Para cada usuario, nos gustaría saber qué grupos de páginas visita, con qué frecuencia se conecta, si tiene un patrón definido de navegación, si este patrón es parecido al de otros usuarios... Realmente no hace falta mucha imaginación para pensar posibles aplicaciones de un análisis de estos resultados.

El comportamiento de un usuario puede definirse con muy diversos niveles de detalle. Nosotros vamos a considerar simplemente que el comportamiento durante una sesión es el conjunto de páginas web que solicita.

Escribe un programa map-reduce que, a partir de unos ficheros log, devuelva los comportamientos de cada usuario (el conjunto de páginas visitadas) y las veces que cada comportamiento se ha repetido.

3 Comportamientos parecidos entre usuarios

Puede ser interesante saber qué usuarios tienen comportamientos parecidos.

Escribe un programa map-reduce que, a partir de unos ficheros log, devuelva para cada comportamiento detectado, la lista de usuarios que han seguido ese comportamiento en una sesión.

4 Integración

Por último, queremos integrar los datos que hemos calculado de forma independiente en los apartados anteriores.

Escribe un programa map-reduce que, a partir de unos ficheros log, devuelva las tuplas:

`user, t_s, behaviour, n_b, user_list`

- `user` es la ip de un usuario
- `t_s` es el número total de sesiones de dicho usuario
- `behaviour` es uno de los posibles comportamientos del usuario (conjunto de páginas)
- `n_b` es el número de veces que el usuario `user` repite el comportamiento `behaviour`
- `user_list` es la lista de usuarios que también han seguido en alguna ocasión el comportamiento `behaviour`.