

BashQL

Manual de Usuario

Esteban Rodríguez Betancourt - B15512

Marie Barquero Rojas - B00799

Lunes 8 de diciembre del 2014

Resumen

BashQL es un conjunto de herramientas que permiten hacer consultas sobre archivos CSV. Cada herramienta está diseñada para que tenga una sola función específica, y al usarlas en conjunto se pueden lograr consultas más complejas.

Las herramientas están programadas en Google Go, y pueden ser utilizadas en cualquier sistema tipo UNIX, como Linux, BSD o OS X.

1. Instalación

La forma más sencilla de instalar BashQL es utilizando el gestor de paquetes de Google Go. Esto además permite recompilar las herramientas para que incorporen futuras mejoras en el compilador de Go. También es posible copiar los ejecutables precompilados para la plataforma correcta: BashQL no tiene dependencias dinámicas aparte de las bibliotecas específicas del sistema.

1.1. Instalación de Google Go

Las siguientes instrucciones son para Go 1.3.3. Es posible que las instrucciones cambien para futuras versiones. Se recomienda revisar el sitio web <http://golang.org/> para tener la información más actualizada.

1.1.1. Instalación precompilada

1. Descargar la versión correspondiente a su sistema operativo y arquitectura desde <https://golang.org/dl/>.
2. Descomprimir el archivo (puede ser en \$HOME).
3. Crear una carpeta aparte para los paquetes de Go (ejm: \$HOME/go-packs).
4. Modificar las variables de entorno para agregar las herramientas de Go:

- a) Agregar en ~/.profile o ~/.bashrc las siguientes líneas, al final. Modificar según sea el caso.

```
1 export GOROOT=$HOME/go
2 export GOPACKS=$HOME/gopacks
3 export PATH=$PATH:$GOROOT/bin:$GOPACKS/bin
```

- b) Cargar las variables de entorno en la sesión actual:

```
1 $ . ~/.profile
```

5. En este punto al correr el comando «go» se deberían desplegar las instrucciones para correr el programa.

1.1.2. Instalación desde código fuente

1. Asegurarse de tener instalado un compilador de C (se recomienda gcc) y mercurial.
2. Clonar el repositorio:

```
1 $ hg clone -u release https://code.google.com/p/go
```

3. Construir la distribución de Go:

```
1 $ cd go/src
2 $ ./all.bash
```

4. Seguir las instrucciones para actualizar las variables de entorno, tal como se indica en los pasos del 3 al 5 de Instalación precompilada.

1.2. Instalación de BashQL

Para descargar BashQL es necesario tener GIT y Go instalados. Luego solamente hay que correr los comandos:

```
1 $ go get github.com/estebarb/bashql/...
2 $ go install github.com/estebarb/bashql/...
```

Con esto se instalan las herramientas que forman parte de BashQL.

2. Uso de BashQL

BashQL está formado por varios comandos, que cumplen con funciones muy específicas. Prácticamente todos los comandos utilizan la entrada y la salida estándar, como método de comunicación entre las herramientas, por lo que es muy sencillo crear nuevos componentes interoperables con BashQL.

2.1. bqlfrom

Nombre bqlfrom - Lee un archivo CSV, lo normaliza y lo escribe en stdout.

Sinopsis bqlfrom [opciones] archivo

Descripción Lee un archivo CSV según el formato especificado en las opciones (o bien el formato estandar), lo transforma al formato estandar y lo escribe en stdout.

- d Indica el caracter que delimita las columnas del archivo CSV. Ejm: «-d ';'» o «-d=';'». Por defecto se utiliza una coma («,») como separador.
- c Indica el caracter de inicio de un comentario. Las líneas comentadas son ignoradas por el lector de CSV. Ejm: «-c '%'» o «-c='%'». Por defecto los CSV no tienen comentarios.

Ejemplo:

```
1 $ cat personal.csv
2 % Personal de la tienda
3 id;nombre;edad;puesto
4 1;Antonio;25;Cajas
5 2;María;26;Cajas
6 $ bqlfrom -d ';' -c='%' personal.csv
7 id,nombre,edad,puesto
8 1,Antonio,25,Cajas
9 2,María,26,Cajas
```

2.2. bqlselect

Nombre bqlselect - Selecciona columnas específicas de un archivo CSV.

Sinopsis bqlselect columnas...

Descripción Recibe por stdin un archivo CSV normalizado y devuelve por stdout los datos de las columnas seleccionadas.

Ejemplo:

```
1 $ bqlfrom -d ';' -c='%' personal.csv \
2 | bqlselect nombre puesto
3 nombre,puesto
4 Antonio,Cajas
5 María,Cajas
```

2.3. bqlwhere

Nombre bqlwhere - Selecciona las filas que cumplen con los criterios dados.

Sinopsis bqlwhere {columna operador argumento}...

Descripción Recibe por stdin un archivo CSV normalizado y devuelve por stdout las filas que cumplen con TODOS los criterios solicitados. Los argumentos van en trios, donde se indica la columna, el operador y el argumento del operador (valor o columna).

Operadores: El comando acepta los siguientes operadores:

- = Igualdad entre el valor de la columna y el valor pasado como argumento (número o cadena).
- < El valor de la columna es menor que el valor pasado como argumento (número o cadena). El texto se compara por orden alfabético.
- > El valor de la columna es mayor que el valor pasado como argumento (número o cadena). El texto se compara por orden alfabético.
- != Desigualdad entre el valor de la columna y el valor pasado como argumento (número o cadena).
- <= El valor de la columna es menor o igual que el valor pasado como argumento (número o cadena). El texto se compara por orden alfabético.
- >= El valor de la columna es mayor o igual que el valor pasado como argumento (número o cadena). El texto se compara por orden alfabético.
- c= Igualdad entre el valor de la columna y el valor en la columna pasada como argumento (número o cadena).
- c< El valor de la columna es menor que el valor en la columna pasada como argumento (número o cadena). El texto se compara por orden alfabético.
- c> El valor de la columna es mayor que el valor en la columna pasada como argumento (número o cadena). El texto se compara por orden alfabético.
- c!= Desigualdad entre el valor de la columna y el valor en la columna pasada como argumento (número o cadena).
- c<= El valor de la columna es menor o igual que el valor en la columna pasada como argumento (número o cadena). El texto se compara por orden alfabético.
- c>= El valor de la columna es mayor o igual que el valor en la columna pasada como argumento (número o cadena). El texto se compara por orden alfabético.

like El texto en la columna hace match contra una expresión regular pasada como argumento.

unlike El texto en la columna NO hace match contra una expresión regular pasada como argumento.

Ejemplo:

```
1 $ bqlfrom poblacion.csv \  
2 | bqlwhere Tipo = 'Provincia' \  
3 Total '<' 500000 \  
4 Nombre '<=' 'Limón'  
5 Nombre,Total,Hombre,Mujer,Tipo  
6 Cartago,491425,241121,250304,Provincia  
7 Heredia,433975,211417,222558,Provincia  
8 Limón,386954,193673,193281,Provincia  
9 Guanacaste,326821,161932,164889,Provincia
```

2.4. bqlwhenever

Nombre bqlwhenever - Selecciona las filas que cumplen con alguno de los criterios dado.

Sinopsis bqlwhenever {columna operador argumento}...

Descripción Recibe por stdin un archivo CSV normalizado y devuelve por stdout las filas que cumplen con AL MENOS UNO de los criterios solicitados. Los argumentos van en trios, donde se indica la columna, el operador y el argumento del operador (valor o columna).

Operadores: El comando acepta los siguientes operadores:

- = Igualdad entre el valor de la columna y el valor pasado como argumento (número o cadena).
- < El valor de la columna es menor que el valor pasado como argumento (número o cadena). El texto se compara por orden alfabético.
- > El valor de la columna es mayor que el valor pasado como argumento (número o cadena). El texto se compara por orden alfabético.
- != Desigualdad entre el valor de la columna y el valor pasado como argumento (número o cadena).
- <= El valor de la columna es menor o igual que el valor pasado como argumento (número o cadena). El texto se compara por orden alfabético.
- >= El valor de la columna es mayor o igual que el valor pasado como argumento (número o cadena). El texto se compara por orden alfabético.

- c=** Igualdad entre el valor de la columna y el valor en la columna pasada como argumento (número o cadena).
- c<** El valor de la columna es menor que el valor en la columna pasada como argumento (número o cadena). El texto se compara por orden alfabético.
- c>** El valor de la columna es mayor que el valor en la columna pasada como argumento (número o cadena). El texto se compara por orden alfabético.
- c!=** Desigualdad entre el valor de la columna y el valor en la columna pasada como argumento (número o cadena).
- c<=** El valor de la columna es menor o igual que el valor en la columna pasada como argumento (número o cadena). El texto se compara por orden alfabético.
- c>=** El valor de la columna es mayor o igual que el valor en la columna pasada como argumento (número o cadena). El texto se compara por orden alfabético.
- like** El texto en la columna hace match contra una expresión regular pasada como argumento.
- unlike** El texto en la columna NO hace match contra una expresión regular pasada como argumento.

Ejemplo:

```

1 $ bqlfrom poblacion.csv \
2 | bqlwhere Tipo = Provincia \
3 | bqlwhenever Nombre '<' 'Cartago' Nombre '>' 'Cartago' \
4 | bqlselect Nombre
5 Nombre
6 San José
7 Alajuela
8 Heredia
9 Puntarenas
10 Limón
11 Guanacaste

```

2.5. bqljoin

Nombre bqljoin - Toma dos archivos CSV y hace un join entre ellos.

Sinopsis bqljoin [opciones] csv1 columna1 csv2 columna2
 bqljoin [opciones] columnaSTDIN csv2 columna2

Descripción Lee dos archivos CSV y hace un join entre ellos según la columna especificada por archivo. Por defecto se hace un join interno, pero con los parámetros se puede modificar para que realice un left join, right join, inner join o outer join / full join. Escribe la salida en stdout.

- d1** Indica el delimitador del primer archivo. Ejm: «-d1=';'. Por defecto es una coma.
- d2** Indica el delimitador del segundo archivo. Ejm: «-d2=';'. Por defecto es una coma.
- d** Indica el delimitador de la entrada estandar. Ejm: «-d=';'. Por defecto es una coma.
- type** Indica el tipo de join que se realizará. El tipo puede ser «inner», «left», «right», «outer» o «full».

Ejemplo:

```

1 $ cat english_numbers.csv
2 number,value
3 zero,0
4 one,1
5 two,2
6 three,3
7 four,4
8 five,5
9 ten,10
10 "one hundred",100
11 "two thousand",2000
12 million,1000000
13 $ cat spanish_numbers.csv
14 numero,value
15 cero,0
16 cuatro,4
17 cinco,5
18 seis,6
19 siete,7
20 ocho,8
21 nueve,9
22 cien,100
23 mil,1000
24 millón,1000000
25 $ bqljoin english_numbers.csv value spanish_numbers.csv value
26 number,value,numero,value
27 zero,0,cero,0
28 one hundred,100,cien,100
29 million,1000000,millón,1000000
30 four,4,cuatro,4
31 five,5,cinco,5
32 $ bqljoin -type left english_numbers.csv value spanish_numbers.csv
    value
33 number,value,numero,value
34 zero,0,cero,0
35 one,1,"",""
36 ten,10,"",""
37 one hundred,100,cien,100
38 million,1000000,millón,1000000
39 two,2,"",""
40 two thousand,2000,"",""
41 three,3,"",""
42 four,4,cuatro,4

```

```

43 five,5,cinco,5
44 $ bqljoin -type right english_numbers.csv value spanish_numbers.csv
    value
45 number,value,numero,value
46 zero,0,cero,0
47 one hundred,100,cien,100
48 "", "",mil,1000
49 million,1000000,millón,1000000
50 four,4,cuatro,4
51 five,5,cinco,5
52 "", "",seis,6
53 "", "",siete,7
54 "", "",ocho,8
55 "", "",nueve,9
56 $ bqljoin -type full english_numbers.csv value spanish_numbers.csv
    value
57 number,value,numero,value
58 zero,0,cero,0
59 one,1,"", ""
60 ten,10,"", ""
61 one hundred,100,cien,100
62 "", "",mil,1000
63 million,1000000,millón,1000000
64 two,2,"", ""
65 two thousand,2000,"", ""
66 three,3,"", ""
67 four,4,cuatro,4
68 five,5,cinco,5
69 "", "",seis,6
70 "", "",siete,7
71 "", "",ocho,8
72 "", "",nueve,9

```

2.6. bqlgroupby

Nombre bqlgroupby - Agrupa una tabla según las columnas especificadas, usando los acumuladores especificados para las demás columnas.

Sinopsis bqlgroupby {-g columnaGrupo}... {-c columnaReducir -f funcionReduccion}

Descripción Lee una tabla por stdin y la agrupa usando las columnas seleccionadas. En las columnas seleccionadas para reducción aplica una función de reducción sobre los valores de dichas columnas.

- g Indica según qué columna se hará el agrupamiento. Ejm: «-g year -g month».
- c Indica de qué columna se tomarán los datos para la función de reducción asociada.
- f Indica qué comando se va a correr como función de reducción. Se puede utilizar cualquier comando que tome sus entradas por stdin (una

por línea) y devuelva el resultado por stdout cuando reciba EOF.
BashQL contiene las siguientes funciones de reducción:

bqlsum Suma los valores recibidos.

bqlcount Cuenta los valores recibidos.

bqldistinct Cuenta los valores únicos recibidos.

bqlavg Calcula el promedio de los valores recibidos.

bqlmax Retorna el mayor valor recibido.

bqlmin Retorna el menor valor recibido.

Ejemplo:

```
1 $ head visitas.csv
2 year,month,day,visits
3 2014,12,06,4253
4 2014,12,05,26286
5 2014,12,04,11716
6 2014,12,03,10206
7 2014,12,02,26719
8 2014,12,01,20816
9 2014,11,30,24176
10 2014,11,29,17463
11 2014,11,28,9194
12 $ bqlfrom visitas.csv \
13 | bqlgroupby -g month \
14 -c visits -f bqlsum \
15 -c visits -f bqlavg \
16 -c visits -f bqlmin \
17 -c visits -f bqlmax
18 month,visits,visits,visits,visits
19 01,6.50238e+06,16134.937965260546,62,32765
20 02,5.976605e+06,16285.027247956403,37,32703
21 03,6.772536e+06,16722.311111111111,8,32713
22 04,6.802518e+06,16196.471428571429,132,32722
23 05,7.27781e+06,16769.14746543779,22,32717
24 06,6.784116e+06,16152.657142857142,78,32581
25 07,7.227396e+06,16652.986175115206,116,32626
26 08,7.131601e+06,16432.260368663596,53,32675
27 09,6.857604e+06,16327.628571428571,30,32419
28 10,6.989063e+06,16103.831797235023,67,32767
29 11,6.852773e+06,16316.12619047619,80,32708
30 12,6.672185e+06,16313.41075794621,21,32743
```

2.7. bqlheader

Nombre bqlheader - Cambia el nombre de las columnas.

Sinopsis bqlheader [opciones] nombres...

Descripción Cambia o agrega el nombre de las columnas del CSV pasado por stdin y lo devuelve por stdout.

-i Sinónimo de insert.

-insert Inserta el encabezado, se usa si el archivo CSV no trae encabezado.

Ejemplo:

```
1 $ bqlfrom visitas.csv \  
2 | bqlgroupby -g month \  
3 -c visits -f bqlsum \  
4 -c visits -f bqlavg \  
5 -c visits -f bqlmin \  
6 -c visits -f bqlmax \  
7 | bqlheader mes total promedio_diario minimo maximo  
8 mes,total,promedio_diario,minimo,maximo  
9 01,6.50238e+06,16134.937965260546,62,32765  
10 02,5.976605e+06,16285.027247956403,37,32703  
11 03,6.772536e+06,16722.311111111111,8,32713  
12 04,6.802518e+06,16196.471428571429,132,32722  
13 05,7.27781e+06,16769.14746543779,22,32717  
14 06,6.784116e+06,16152.657142857142,78,32581  
15 07,7.227396e+06,16652.986175115206,116,32626  
16 08,7.131601e+06,16432.260368663596,53,32675  
17 09,6.857604e+06,16327.628571428571,30,32419  
18 10,6.989063e+06,16103.831797235023,67,32767  
19 11,6.852773e+06,16316.12619047619,80,32708  
20 12,6.672185e+06,16313.41075794621,21,32743
```

2.8. bqlsort

Nombre bqlsort - Ordena las filas según las columnas pasadas como argumento.

Sinopsis bqlsort [opciones] columnas...

Descripción Recibe por stdin un archivo CSV normalizado y devuelve por stdout el archivo con las filas ordenadas de forma ascendente según las columnas pasadas como parámetro.

Opciones: El comando acepta las siguientes opciones:

-r Ordena de forma descendente

-reverse Sinónimo de -r.

Ejemplo:

```
1 $ bqlfrom poblacion.csv \  
2 | bqlwhere Tipo = Provincia \  
3 | bqlwhenever Nombre '<' 'Cartago' \  
4 | bqlwhenever Nombre '>' 'Cartago' \  
5 | bqlsort Nombre  
6 Nombre>Total,Hombre,Mujer,Tipo
```

7	Alajuela,847660,420636,427024,Provincia
8	Guanacaste,326821,161932,164889,Provincia
9	Heredia,433975,211417,222558,Provincia
10	Limón,386954,193673,193281,Provincia
11	Puntarenas,410914,205975,204939,Provincia
12	San José,1403963,671434,732529,Provincia