# Find the most popular status entry (awk, sort, cat, csvcut, head)

To do this analysis efficiently, we'll use the command line language called `awk`, a tool that allows you to filter, extract and transform data files. awk is a very useful tool to put in your bag of tricks. But let's watch the following video lecture first!



Find the most popular status entry

To start, let's look at a very simple awk program to output every line of our
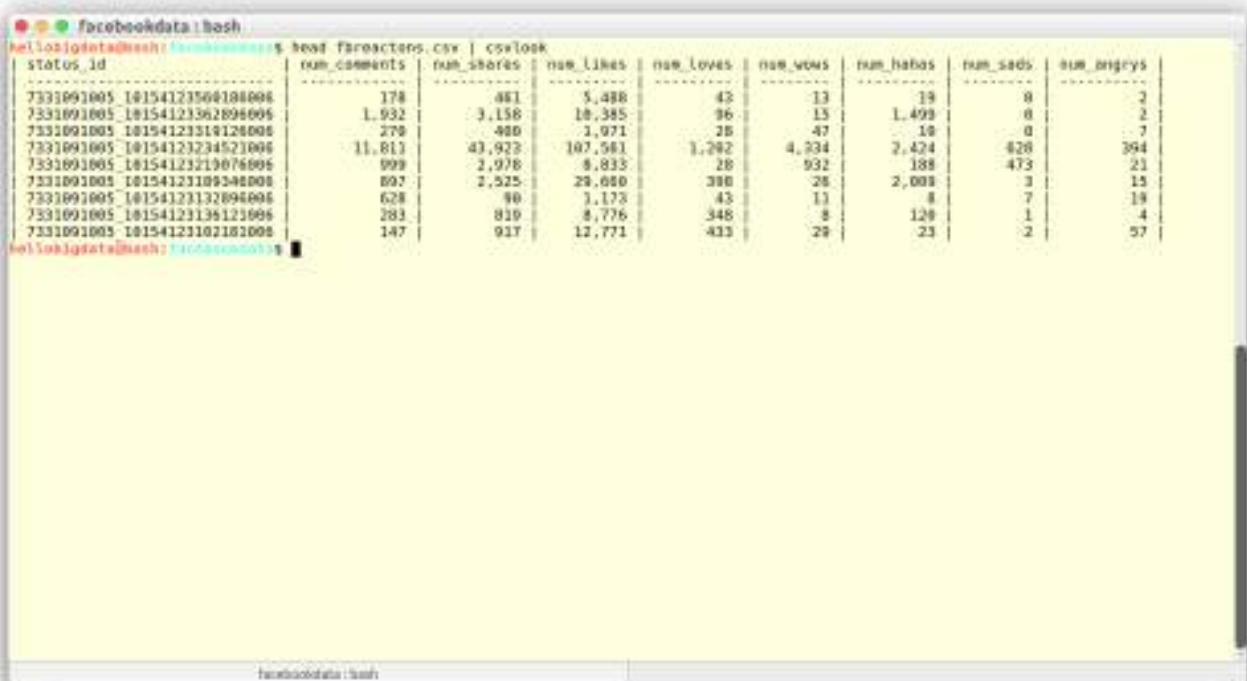
-F option:

```
awk -F "," '{ print; }' facebookdata.csv
```

You should see the entire file being output to the screen. To only output the status ids (column 1), use the dollar sign ( `$` ) to denote columns as follows:

```
awk -F "," '{ print $1; }' facebookdata.csv | head
```

However, since the dataset has quoted ( `"text"` ) cells we will use `csvcut` to extract the columns, e.g., we want to extract the column `1,8-15` into a file called `fbreactions.csv` . The idea is to sum-up all the reactions (columns `8 + ... + 15` ) on each FB status and then find the status which had the maximum number of reactions.

```
csvcut -c 2,8-15 facebookdata.csv > fbreactons.csv
```



Extract all the reactions into a file fbreactions.csv

To calculate the total number of reactions on each entry (status), all we need to do is horizotally add up all the numbers from the columns #8-15 and we do this easily with awk, as follows:

```
awk -F "," '{ total = total + $2 + $3 + $4 + $5 + $6 +$7 +$8 +$9; print $1" " total; total=0
```

```
fbreactons.csv | \
head
```

Let's pay attention to the `awk` statetment, which not only sums up the columns side by side, but also on each line prints two output ( `status id` and `total` number of reaction on that row). Finally, at the end of each iteration, it nulls the `total=0` .

To get the status with `max` reactions, next, we sort the status ids, based on the number of reactions (column 2) using the `sort -n -r -t"," -k 2` function, which tells the system to sort out the piped ( `|` ) output numerically ( `-n` ), on the column 2 ( `-k 2` ) wich is delimited by a commma ( `,` ):

```
awk -F "," '{ total = total + $2 + $3 + $4 + $5 + $6 +$7 +$8 +$9; print $1"," total; total=0
fbreactons.csv | \
sort -n -r  -t"," -k 2 | \
head -n 1
```



The final output, tells us that the status id: `7331091005_10154089857531006` had the maximum number of reaction of total `668121` .

If we now use `grep`, we can easily find the message which had the largest number of reactions.

```
cat facebookdata.csv | grep 7331091005_10154089857531006
```

Let's make it little more beautiful using `csvcut`:



FB Awk total find

However, we want to make it more interesting! let's efficently pipe all the steps shown above into a single command and find the message as follows:

```
cat facebookdata.csv | \
csvcut -c 2,8-15 | \
awk -F "," '{ total = total + $2 + $3 + $4 + $5 + $6 +$7 +$8 +$9; print $1","total; total=0 }
sort -n -r  -t"," -k 2 | \
head -n 1
```

# Do you want to know more? #

📎 'awk' man page ⬇ ↱