Finding the top most crime per city (csvcut, sort, tail, awk)

Our next compute quest is to find the total number of crimes, given a city name. Let's watch the following video lecture first!

Video lecture: Finding the top most crime per city

This essentially will require us to compute the summation of a column which contains the city names. Ideally, we would like to have a function like top_crime () that would take an argument (e.g., Sydney) and tell us the top crime for that city. Let's start with a basic version first:

```
cat crimedata-au.csv | \
csvcut -c "Case Incident Type","Sydney" | \
sort -n -t "," -k 2
```





```
🔵 🔵 🔵 afpdata : bash
hellobigdata@bash:
                     pdata$ cat crimedata-au.csv | csvcut -c "Case Incident Type","Sydney" | sort -n -r -t "," -k 2
Drugs ? Imported,191
Money Laundering,83
Fraud,63
Child Sex Offences ? Online Child Sex Exploitation,54
Identity Crime,13
Migration Crime, 10
Intellectual Property,7
Terrorism Domestic,6
Slavery/Human Trafficking,6
Emerging Crime,6
Information And Communications Technology,5
Corruption,5
Offences On Commonwealth Employee,4
Child Sex Offences ? Not Child Sex Tourism,4
War Crimes,2
Threats,2
Terrorism Financing,2
Drugs ? Trafficked,2
Weapons/Prohibited Items In The Aviation Environ.,1
Transnational ? Child Sex Tourism,1
Theft ? General.1
Terrorism International,1
People Smuggling,1
Counterfeit Currency,1
Corporate Or Bankruptcy,1
Weapons ? Trafficked,0
Transnational ? Sexual Servitude,0
Transnational Organised Crime,0
                          afpdata: bash
```

Finding the total number of crimes, for Sydney

Note that the csvcut function can take column titles with the -c option. Now, if we just want the top crime we add tail -n 1 at the end, as follows:

```
cat crimedata-au.csv | \
csvcut -c "Case Incident Type","Sydney" | \
sort -n -t "," -k 2 | \
tail -n 1
```

Now, let's build our desired function. For knowing more about bash function, please refer to the tutorials section.

```
function topcrime() { cat $1 | \
   csvcut -c "Case Incident Type",$2 | \
   sort -n -t "," -k 2 | tail -n 1| \
   awk -F',' '{print $1 "= " $2}'; }

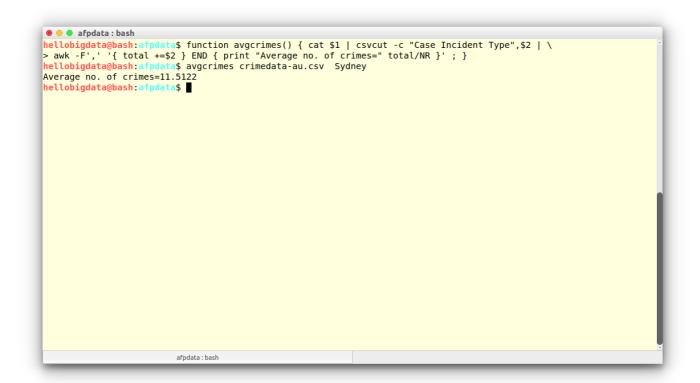
topcrime crimedata-au.csv Sydney
```

arguements were automatically placed into the \$1 and \$2 variables inside

the function, however indside awk 's call \$1 and \$2 represented the first (crime type) and second (number of crimes) columns of the piped output.

Now, we can slightly modify the above function so that it can also compute, given any city name, the average number of crimes.

```
function avgcrimes() { cat $1 | \
  csvcut -c "Case Incident Type",$2 | \
  awk -F',' '{ total +=$2 } END { print "Average no. of crimes=" total/NR }'; }
  avgcrimes crimedata-au.csv Sydney
```



Finding the average number of crimes, for Sydney

Here, everything remains as above, except there's no sort function and the awk function computes the sum of crimes (col 2) and then divide the sum by the total number of rows (NR).