

```

1  #!/bin/bash
2
3  # In `hai.csv`, the location is split up into 3 lines,
4  # which is fine, but I just wanted to get rid of them for
5  # readability.
6  echo '- Getting rid of unnecessary newlines on csv (this may take a moment)...'
7
8  # Grab the header and drop it into a new csv, `temp.csv`.
9  head -1 downloads/hai.csv > temp.csv
10 # Skip the header row of `hai.csv`,
11 # join every 3 lines for readability (the address is split into 2 lines),
12 # replace 'Not Available' with 'NA' (for a smaller file size later),
13 # and append to `temp.csv`
14 tail -n +2 downloads/hai.csv | paste -d, - - - | sed -E 's/Not Available/NA/g' >>
    temp.csv
15
16 # Make csv of hospitals
17 echo '- Extracting hospital info (this may take a moment)...'
18 # Extract the following columns from `temp.csv`:
19 # 1: Provider ID
20 # 2: Hospital Name
21 # 3: Address
22 # 4: City
23 # 5: State
24 # 6: ZIP Code
25 # 16: Location
26 # Then sort the csv by Provider ID,
27 # filter for uniques,
28 # and create a new csv, `hospitals_temp.csv`.
29 csvcut -c 1,2,3,4,5,6,16 temp.csv | csvsort -c 1 | uniq > hospitals_temp.csv
30
31 # Make csv of footnotes
32 echo '- Filtering rows with footnotes... '
33 # Grep for the central line data that we care about (skipping confidence intervals),
34 # cut out columns 1 (Provider ID) and 13 (Footnotes),
35 # grep Footnotes column (which is now column 2) for cells with numbers (footnotes),
36 # sort Footnotes column,
37 # filter for uniques,
38 # reverse grep for footnote about confidence intervals because we don't need that,
39 # and create a new csv, `footnotes_temp.csv`.
40 csvgrep -c 10 -r "HAI_1_DOPC_DAYS|HAI_1_NUMERATOR|HAI_1_SIR" temp.csv | csvcut -c 1,13 |
    csvgrep -c 2 -r '\d+' | csvsort -c 2 | uniq | csvgrep -c 2 -i -m "8 - The lower limit of
    the confidence interval cannot be calculated if the number of observed infections equals
    zero." > footnotes_temp.csv
41 # (Note to self: Once you filter out the measures, there's no `3, 8` combo footnote.)
42
43 # Return a List of footnotes
44 echo '- Creating a list of footnotes (this may take a moment)...'
45 # Extract the Footnotes column in `footnotes_temp.csv`
46 # but not the header row,
47 # sort for uniques,
48 # and create `footnotes_list.txt`
49 csvcut -c 2 footnotes_temp.csv | tail -n +2 | sort -u > footnotes_list.txt
50
51 # Make the data a bit more readable by cutting out the footnotes' long explanations.
52 # We'll place the results in `footnotes.csv`.

```

```

53 echo '- Shortening footnotes...'
54 sed -E 's/12 - This measure does not apply to this hospital for this reporting
period./12/g' footnotes_temp.csv | sed -E 's/13 - Results cannot be calculated for this
reporting period./13/g' | sed -E 's/3 - Results are based on a shorter time period than
required./3/g' | sed -E 's/5 - Results are not available for this reporting period./5/g'
| sed -E 's/8 - The lower limit of the confidence interval cannot be calculated if the
number of observed infections equals zero./8/g' | sed -E 's/, /,/g' > footnotes.csv
55
56
57 echo '- Formatting latitude and longitude...'
58 # The hospital data comes with geographic coordinates for each hospital,
59 # formatted like `(31.21537937900007, -85.36146587999997)`.
60 # I very likely did not need a separate Python script for this,
61 # and may not have needed to split the coordinates at all.
62 python location.py
63 # `location.py` outputs `hospitals_info.csv`.
64
65 echo '- Filtering hospital type info...'
66 # `hospital_general.csv` has the same issue as `hai.csv` in that
67 # the location has a couple newlines in it. Not going to join lines
68 # for readability here since the following will extract columns
69 # without newlines.
70 # Extract the following columns out of `downloads/hospital_general.csv`:
71 # 1: Provider ID
72 # 9: Hospital Type
73 # 10: Hospital Ownership
74 # 11: Emergency Services
75 # and put them in `hospitals_type.csv`.
76 csvcut -c 1,9,10,11 downloads/hospital_general.csv > hospitals_type.csv
77
78 echo '- Joining hospital general information data...'
79 # Left join `hospitals_info.csv` and `hospitals_type.csv` on their first columns
  (Provider ID),
80 # remove the 9th column because that's a dupe of Provider ID,
81 # and put it all in `hospitals.csv`.
82 csvjoin -c 1,1 hospitals_info.csv hospitals_type.csv | csvcut -C 9 > hospitals.csv
83 # Note: There are more rows in `hospitals_type.csv` than `hospitals_info.csv`
84 # because not all hospitals report central line measures to this particular agency.
85
86 echo '- Filtering out CLABSI measures (this may take a moment)...'
87 # Grep for rows with one of the central line measures in the 10th column,
88 # then extract the following columns:
89 # 1: Provider ID
90 # 10: Measure ID
91 # 12: Score
92 # and dump that all in `clabsi_temp.csv`
93 csvgrep -c 10 -r "HAI_1_DOPC_DAYS|HAI_1_NUMERATOR|HAI_1_SIR" temp.csv | csvcut -c 1,10,12
  > clabsi_temp.csv
94
95 echo '- Filtering CLABSI SIR...'
96 # Grep for the SIR rows,
97 # remove the label column,
98 # and dump to `clabsi_sir.csv`.
99 csvgrep -c 2 -m "HAI_1_SIR" clabsi_temp.csv | csvcut -C 2 > clabsi_sir.csv
100
101 echo '- Filtering CLABSI days...'
102 # Grep for the days rows,

```

```
103 # remove the label column,  
104 # and dump to `clabsi_days.csv`.  
105 csvgrep -c 2 -m "HAI_1_DOPC_DAYS" clabsi_temp.csv | csvcut -C 2 > clabsi_days.csv  
106  
107 echo '- Filtering CLABSI observed cases...'  
108 # Grep for the observed rows,  
109 # remove the label column,  
110 # and dump to `clabsi_observed.csv`.  
111 csvgrep -c 2 -m "HAI_1_NUMERATOR" clabsi_temp.csv | csvcut -C 2 > clabsi_observed.csv  
112  
113 # Join the data we need into one big table.  
114 echo '- Joining tables...'  
115  
116 # Write a header row to `hospitals_clabsi.csv`.  
117 echo  
118 'provider_id,hospital_name,street,city,state,zip_code,lat,lng,type,ownership,emergency_se  
119 rvices,observed,days,sir,footnotes' > hospitals_clabsi.csv  
120  
121 # Left join on the "Provider ID" column: `hospitals.csv`, the csvs with all the CLABSI  
122 scores, and the footnotes,  
123 # remove the columns that are dupes of "Provider ID" after the join,  
124 # and stream everything (but the header row) into `hospitals_clabsi.csv`.  
125 csvjoin -c "Provider ID" --left hospitals.csv clabsi_observed.csv clabsi_days.csv  
126 clabsi_sir.csv footnotes.csv | csvcut -C 12,14,16,18 | tail -n +2 >> hospitals_clabsi.csv  
127  
128 echo '- Cleaning up data...'  
129  
130 # Does everything look good?  
131 csvclean hospitals_clabsi.csv  
132  
133 # It does.  
134 mv hospitals_clabsi_out.csv hospitals_clabsi.csv  
135  
136 echo '- Removing unnecessary files...'  
137 rm clabsi_days.csv clabsi_observed.csv clabsi_sir.csv clabsi_temp.csv hospitals.csv  
138 hospitals_info.csv hospitals_temp.csv hospitals_type.csv footnotes_temp.csv footnotes.csv  
139 temp.csv  
140  
141 echo 'Data processing complete. Check `hospitals_clabsi.csv` for complete table, and  
142 `foontotes_list.txt` for a list of what the footnotes are.'
```