7/18/2019 build.sh

```
#!/bin/bash
 1
 2
 3
   # In `hai.csv`, the location is split up into 3 lines,
   # which is fine, but I just wanted to get rid of them for
 5
   # readability.
   echo '- Getting rid of unnecessary newlines on csv (this may take a moment)...'
 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),
   # replace 'Not Available' with 'NA' (for a smaller file size later),
# 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,
   # and create a new csv, `hospitals_temp.csv`.
28
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...'
   # 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,
   # reverse grep for footnote about confidence intervals because we don't need that,
38
39 # and create a new csv, `footnotes temp.csv`.
   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
   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`
   csvcut -c 2 footnotes_temp.csv | tail -n +2 | sort -u > footnotes_list.txt
49
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`.
```

7/18/2019 build.sh

```
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
    # `location.py` outputs `hospitals info.csv`.
63
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
    echo '- Joining hospital general information data...'
 78
 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
    echo '- Filtering out CLABSI measures (this may take a moment)...'
 86
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`
    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,
    # and dump to `clabsi sir.csv`.
    csvgrep -c 2 -m "HAI 1 SIR" clabsi temp.csv | csvcut -C 2 > clabsi sir.csv
99
100
101 echo '- Filtering CLABSI days...'
    # Grep for the days rows,
```

7/18/2019 build.sh

```
103 # remove the label column,
    # and dump to `clabsi days.csv`.
    csvgrep -c 2 -m "HAI 1 DOPC DAYS" clabsi temp.csv | csvcut -C 2 > clabsi days.csv
106
107
    echo '- Filtering CLABSI observed cases...'
    # Grep for the observed rows,
108
109
    # remove the label column,
    # and dump to `clabsi observed.csv`.
110
    csvgrep -c 2 -m "HAI_1_NUMERATOR" clabsi_temp.csv | csvcut -C 2 > clabsi_observed.csv
111
112
113
    # Join the data we need into one big table.
    echo '- Joining tables...'
114
115
116
    # Write a header row to `hospitals clabsi.csv`.
117
     'provider id, hospital name, street, city, state, zip code, lat, lng, type, ownership, emergency se
     rvices,observed,days,sir,footnotes' > hospitals clabsi.csv
118
    # Left join on the "Provider ID" column: `hospitals.csv`, the csvs with all the CLABSI
119
     scores, and the footnotes,
120
    # remove the columns that are dupes of "Provider ID" after the join,
121 # and stream everything (but the header row) into `hospitals clabsi.csv`.
    csvjoin -c "Provider ID" --left hospitals.csv clabsi observed.csv clabsi days.csv
122
    clabsi_sir.csv footnotes.csv | csvcut -C 12,14,16,18 | tail -n +2 >> hospitals_clabsi.csv
123
124
    echo '- Cleaning up data...'
125
126
    # Does everything Look good?
    csvclean hospitals clabsi.csv
127
128
129
    # It does.
130
    mv hospitals clabsi out.csv hospitals clabsi.csv
131
132
    echo '- Removing unnecessary files...'
133
    rm clabsi days.csv clabsi observed.csv clabsi sir.csv clabsi temp.csv hospitals.csv
     hospitals info.csv hospitals temp.csv hospitals type.csv footnotes temp.csv footnotes.csv
     temp.csv
134
     echo 'Data processing complete. Check `hospitals clabsi.csv` for complete table, and
135
     `foonotes list.txt` for a list of what the footnotes are.'
136
```