

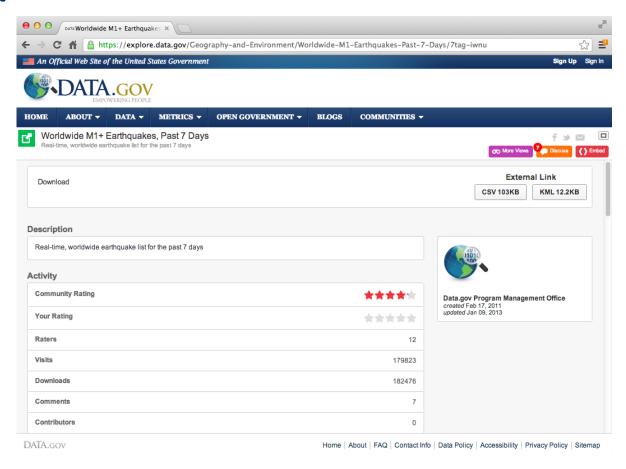
Summarizing data

Jeffrey Leek, Assistant Professor of Biostatistics Johns Hopkins Bloomberg School of Public Health

Why summarize?

- · Data are often too big to look at the whole thing
- The first step in an analysis is to find problems
- · When you do these summaries you should be looking for
 - Missing values
 - Values outside of expected ranges
 - Values that seem to be in the wrong units
 - Mislabled variables/columns
 - Variables that are the wrong class

Earthquake data



http://earthquake.usgs.gov/earthquakes/feed/v1.0/

Earthquake data

```
fileUrl <- "http://earthquake.usgs.gov/earthquakes/feed/v1.0/summary/1.0_week.csv"
download.file(fileUrl,destfile="./data/earthquakeData.csv",method="curl")
dateDownloaded <- date()
dateDownloaded</pre>
```

```
[1] "Wed Aug 28 14:31:30 2013"
```

```
eData <- read.csv("./data/earthquakeData.csv")
```

Looking at data - the whole thing

eData

```
time latitude longitude
                                                     depth mag magType
                                                     5.500 2.20
1
     2013-08-28T18:17:42.700Z
                                 37.405
                                         -121.755
                                                                      Md
2
     2013-08-28T17:54:40.400Z
                                 37.404
                                         -121.756
                                                     5.600 2.30
                                                                      Md
3
     2013-08-28T17:54:14.000Z
                                 38.781
                                         -122.927
                                                     1.700 1.50
                                                                      Md
4
     2013-08-28T17:40:44.860Z
                                 37.025
                                         -117.740
                                                     7.000 1.16
                                                                      ml
5
     2013-08-28T16:36:34.778Z
                                 39.669
                                         -119.678
                                                    13.813 1.70
                                                                      ml
6
     2013-08-28T16:35:17.000Z
                                 61.603
                                         -141.177
                                                     0.100 2.20
                                                                      MΊ
7
                                                                      MΊ
     2013-08-28T16:12:13.300Z
                                 33.851
                                         -117.811
                                                     0.100 1.40
8
     2013-08-28T16:09:26.000Z
                                 62.611
                                         -151.316
                                                    86.800 1.20
                                                                      Ml
9
     2013-08-28T16:03:17.100Z
                                 36.438
                                         -121.004
                                                     0.000 1.80
                                                                      Md
                                         -116.446
                                                    11.500 1.20
                                                                      Ml
10
     2013-08-28T15:21:21.200Z
                                 33.517
11
     2013-08-28T15:13:24.900Z
                                 38.774
                                         -122.715
                                                     2.200 1.90
                                                                      Md
                                                                      MΊ
12
     2013-08-28T14:39:54.000Z
                                 61.599
                                         -141.230
                                                     0.100 2.50
13
     2013-08-28T14:32:21.050Z
                                 39.213
                                            74.585
                                                    55.910 4.50
                                                                      mb
14
     2013-08-28T14:26:48.500Z
                                 38.776
                                         -122.716
                                                     2.400 2.00
                                                                      Md
15
     2013-08-28T13:58:41.300Z
                                 38.750
                                         -122.701
                                                     2.500 1.00
                                                                      Md
16
     2013-08-28T13:25:49.100Z
                                 35.337
                                         -117.914
                                                     8.900 1.50
                                                                      Ml
17
                                 42.028
                                            85.850
                                                    25.860 4.60
     2013-08-28T12:57:40.810Z
                                                                      mb
18
     2013-08-28T12:50:13.860Z
                                 37.441
                                          144.495
                                                    10.110 4.60
                                                                      mb
19
                                                                      MΊ
     2013-08-28T12:47:42.000Z
                                 19.367
                                         -155.027
                                                     5.500 2.20
                                                                                                  5/21
20
     2013-08-28T12:44:51.000Z
                                 59.883
                                         -152.292
                                                    30.900 1.20
                                                                      Ml
```

Looking at data - dim(),names(),nrow(),ncol()

```
dim(eData)
[1] 1076
           14
names (eData)
 [1] "time"
                  "latitude" "longitude" "depth"
                                                        "mag"
                  "nst"
                              "gap"
                                           "dmin"
                                                        "rms"
 [6] "magType"
                  "id"
                              "updated"
                                           "place"
[11] "net"
nrow(eData)
[1] 1076
```

Looking at the data - quantile()

quantile(eData\$latitude)

0% 25% 50% 75% 100% -57.96 33.84 38.82 60.01 67.57

Looking at the data - summary()

summary(eData)

```
time
                                     latitude
                                                     longitude
2013-08-21T18:38:20.100Z:
                                  Min.
                                         :-58.0
                                                  Min.
                                                          :-180
                                  1st Qu.: 33.8
2013-08-21T18:38:57.000Z:
                                                  1st Qu.:-148
2013-08-21T18:55:41.700Z:
                                                  Median :-122
                                  Median: 38.8
2013-08-21T19:06:43.800Z:
                                  Mean
                                         : 41.4
                                                  Mean
                                                          :-114
2013-08-21T19:11:22.100Z:
                                  3rd Ou.: 60.0
                                                  3rd Ou.:-116
2013-08-21T19:11:35.900Z:
                                         : 67.6
                                                          : 180
                                  Max.
                                                  Max.
(Other)
                         :1070
    depth
                                   magType
                                                    nst
                      mag
       : -2.1
Min.
                Min.
                        :1.00
                                M
                                        :633
                                               Min.
                                                       : 0.0
1st Qu.: 3.6
                1st Qu.:1.20
                                               1st Qu.: 12.0
                                Md
                                        :310
                                               Median: 19.0
Median: 9.0
                 Median :1.60
                                        : 80
       : 24.7
                                                       : 27.5
                        :1.91
                                ml
                                        : 38
Mean
                 Mean
                                               Mean
3rd Ou.: 21.4
                 3rd Qu.:2.20
                                               3rd Qu.: 32.0
                                          4
       :592.3
                        :6.10
                                                       :249.0
Max.
                 Max.
                                Mw
                                               Max.
                                 (Other):
                                               NA's
                                                       :519
                      dmin
                                                       net
                                      rms
     qap
Min.
       : 0.0
                        : 0.0
                                Min.
                                        :0.010
                                                         :369
                 Min.
                                                 ak
                 1st Qu.: 0.0
1st Ou.: 64.8
                                1st Qu.:0.120
                                                         :230
Median: 94.0
                 Median: 0.1
                                Median :0.230
                                                         :186
                                                 nc
                                                                                                 8/21
       :117.9
                        : 0.7
                                        :0.344
                                                         : 91
                 Mean
Mean
                                Mean
                                                 us
```

Looking at data - class()

```
class(eData)
```

```
[1] "data.frame"
```

```
sapply(eData[1,],class)
```

```
time latitude longitude depth mag magType nst

"factor" "numeric" "numeric" "numeric" "factor" "integer"

gap dmin rms net id updated place

"numeric" "numeric" "factor" "factor" "factor" "factor"
```

Looking at data - unique(), length(), table()

unique(eData\$net)

[1] nc nn ak ci us hv uw pr se uu nm ld at mb Levels: ak at ci hv ld mb nc nm nn pr se us uu uw

length(unique(eData\$net))

[1] 14

table(eData\$net)

ak at ci hv ld mb nc nm nn pr se us uu uw 369 1 230 27 2 4 186 6 38 71 4 91 12 35

Looking at data - table()

table(eData\$net,eData\$mag)

```
1 1.01 1.03 1.05 1.07 1.1 1.16 1.18 1.2 1.22 1.29 1.3 1.31 1.32
ak 22
                                                                  35
           0
                 0
                       0
                                36
                                                 37
                                                              0
                                                                         0
                                                                               0
at
    0
           0
                0
                       0
                                 0
                                       0
                                                  0
                                                        0
                                                                         0
                                                                               0
ci 25
                                28
                                                 30
                0
                       0
                                                        0
                                                                  25
                                                                         0
                                                                               0
    2
                0
                       0
                            0
                                 0
                                                  0
                                                                         0
                                                                               0
hv
           0
                                       0
                                                        0
                                                                   1
ld
    0
           0
                0
                       0
                                                  0
                                                        0
                                                                         0
                                                                               0
    2
mb
           0
                0
                       0
                            0
                                 0
                                       0
                                             0
                                                  0
                                                        0
                                                                   1
                                                                         0
                                                                               0
nc 24
                                12
                                                 28
                                                                         0
                                                                               0
           0
                       0
                                                                  14
    0
           0
                0
                       0
                            0
                                 0
                                       0
                                                  0
                                                        0
                                                              0
                                                                   1
                                                                         0
                                                                               0
nm
    0
                                 0
                                                  0
                                                                   0
                                                                         1
nn
    0
           0
                0
                       0
                            0
                                 0
                                       0
                                             0
                                                        0
                                                              0
                                                                         0
                                                                               0
pr
    0
                0
                       0
                                       0
                                                        0
                                                                         0
                                                                               0
           0
                                             0
se
    0
                      0
                                 0
                                                                         0
                                                                               0
           0
                                             0
us
                                       0
                                                                         0
                                                                               0
           0
                       0
                                             0
uu
    6
           0
                0
                       0
                            0
                                 3
                                       0
                                             0
                                                              0
                                                                   1
                                                                         0
                                                                               0
uw
              1.45 1.47 1.5 1.52 1.58 1.6 1.63 1.65 1.68 1.7 1.78
       0
          32
                           31
                                            22
                                                                   24
                                                                              22
ak
                                                                               0
                                                         0
at
                                                                                                             11/21
          37
                           27
ci
                  0
                        0
                                  0
                                        0
                                             9
                                                   0
                                                         0
                                                               0
                                                                  13
                                                                          0
```

Looking at data - any(), all()

eData\$latitude[1:10]

[1] 37.40 37.40 38.78 37.02 39.67 61.60 33.85 62.61 36.44 33.52

eData\$latitude[1:10] > 40

[1] FALSE FALSE FALSE FALSE TRUE FALSE TRUE FALSE FALSE

any(eData\$latitude[1:10] > 40)

[1] TRUE

Looking at data - all()

eData\$latitude[1:10] > 40

[1] FALSE FALSE FALSE FALSE TRUE FALSE TRUE FALSE FALSE

all(eData\$latitude[1:10] > 40)

[1] FALSE

Looking at subsets - &

```
eData[eData$latitude > 0 & eData$longitude > 0,c("latitude","longitude")]
```

```
latitude longitude
      39.213
                 74.58
13
      42.028
                 85.85
17
18
      37.441
                144.50
45
      41.827
                139.91
      28.398
107
               99.09
123
      32.739
                 56.41
      28.252
               99.33
135
      37.595
163
                142.05
      44.506
                149.06
181
      47.122
                152.66
205
      2.171
                128.61
208
      56.546
                112.70
247
      34.024
269
               87.94
      33.176
                 94.10
302
316
      30.038
                 97.88
      34.382
                141.01
345
      49.577
355
                155.53
      38.826
375
                69.93
                144.04
431
      22.145
                                                                                               14/21
449
      35.428
                140.16
```

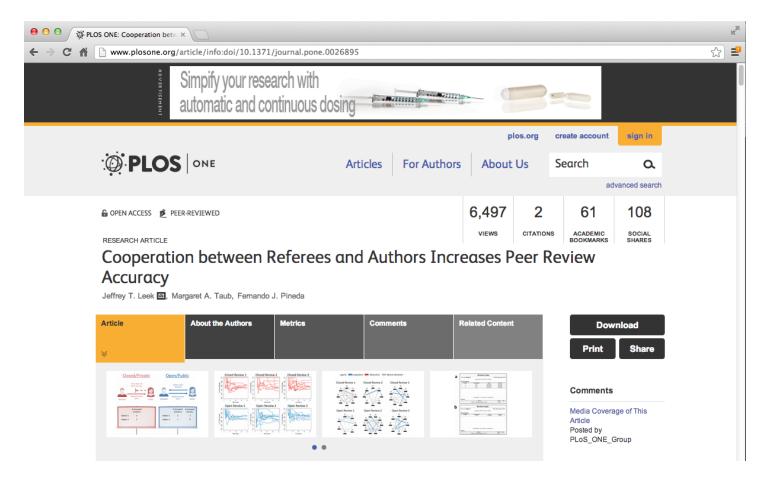
Looking at subsets - |

```
eData[eData$latitude > 0 | eData$longitude > 0,c("latitude","longitude")]
```

```
latitude longitude
       37.405 -121.755
1
       37.404 -121.756
2
       38.781 -122.927
3
       37.025 -117.740
4
5
      39.669 -119.678
       61.603 -141.177
6
       33.851 -117.811
7
       62.611 -151.316
8
      36.438 -121.004
9
      33.517 -116.446
10
       38.774 -122.715
11
12
       61.599 -141.230
13
      39.213
                74.585
      38.776
              -122.716
14
       38.750
              -122.701
15
      35.337 -117.914
16
17
      42.028
                85.850
      37.441
              144.495
18
19
      19.367 -155.027
                                                                                             15/21
20
       59.883 -152.292
```

Peer review experiment data

· Data on submissions/reviews in an experiment



http://www.plosone.org/article/info:doi/10.1371/journal.pone.0026895

Peer review data

```
fileUrl1 <- "https://dl.dropbox.com/u/7710864/data/reviews-apr29.csv"
fileUrl2 <- "https://dl.dropbox.com/u/7710864/data/solutions-apr29.csv"
download.file(fileUrl1,destfile="./data/reviews.csv",method="curl")
download.file(fileUrl2,destfile="./data/solutions.csv",method="curl")
reviews <- read.csv("./data/reviews.csv"); solutions <- read.csv("./data/solutions.csv")
head(reviews,2)</pre>
```

```
X.html.
1 <head><title>Found</title></head>
2 <body>
```

```
head(solutions,2)
```

```
X.html.
1 <head><title>Found</title></head>
2 <body>
```

Find if there are missing values - is.na()

<pre>is.na(reviews\$time_left[1:10])</pre>
logical(0)
<pre>sum(is.na(reviews\$time_left))</pre>
[1] 0
table(is.na(reviews\$time_left))

Important table()/NA issue

```
table(c(0,1,2,3,NA,3,3,2,2,3))
```

```
0 1 2 3
1 1 3 4
```

```
table(c(0,1,2,3,NA,3,3,2,2,3),useNA="ifany")
```

```
0 1 2 3 NA>
1 1 3 4 1
```

Summarizing columns/rows - rowSums(),rowMeans(),colSums(),colMeans()

· Important parameters: x, na.rm

colSums(reviews)

Error: 'x' must be numeric

Summarizing columns/rows - rowSums(),rowMeans(),colSums(),colMeans()

colMeans(reviews,na.rm=TRUE)

Error: 'x' must be numeric

rowMeans(reviews,na.rm=TRUE)

Error: 'x' must be numeric