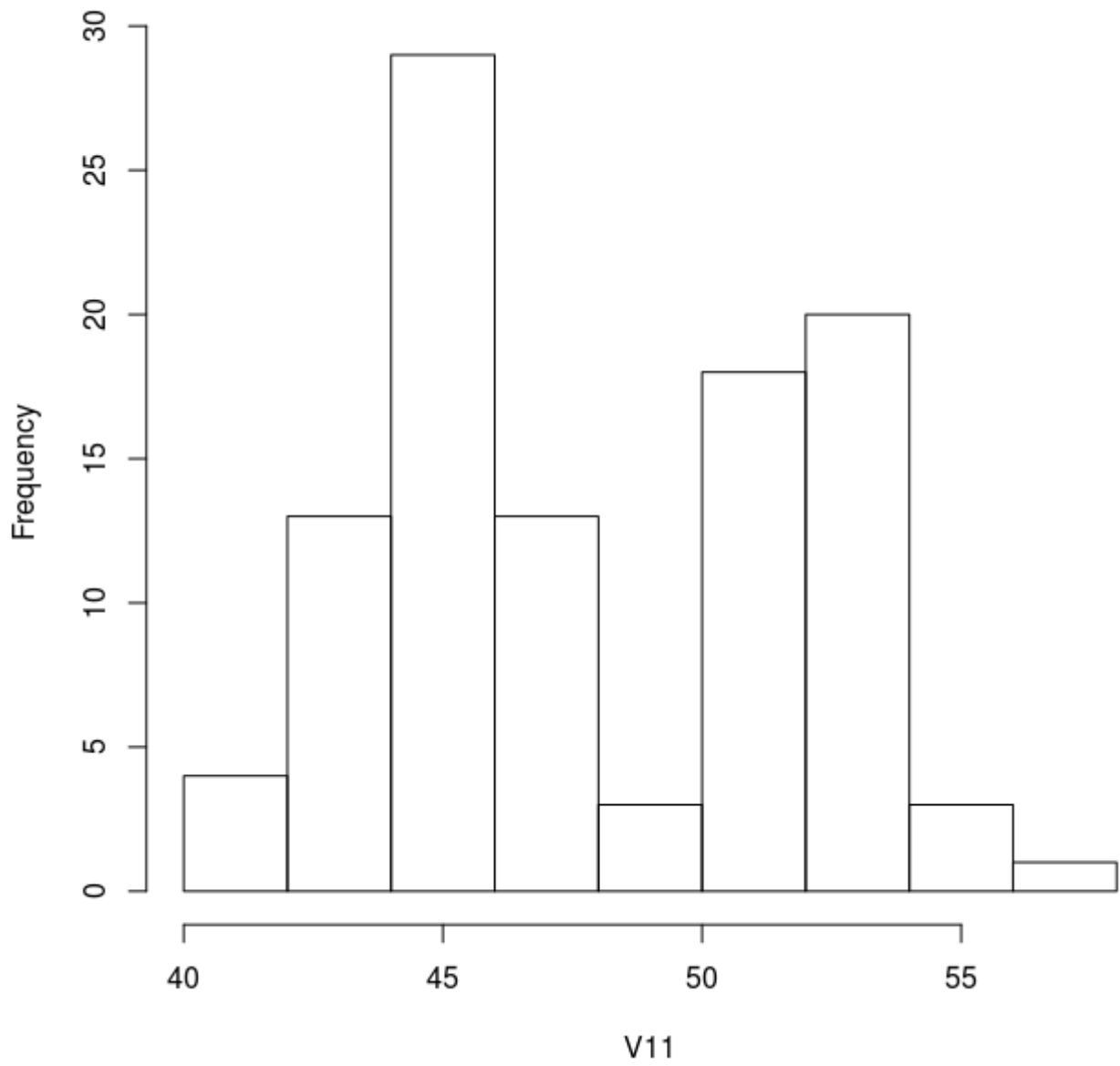
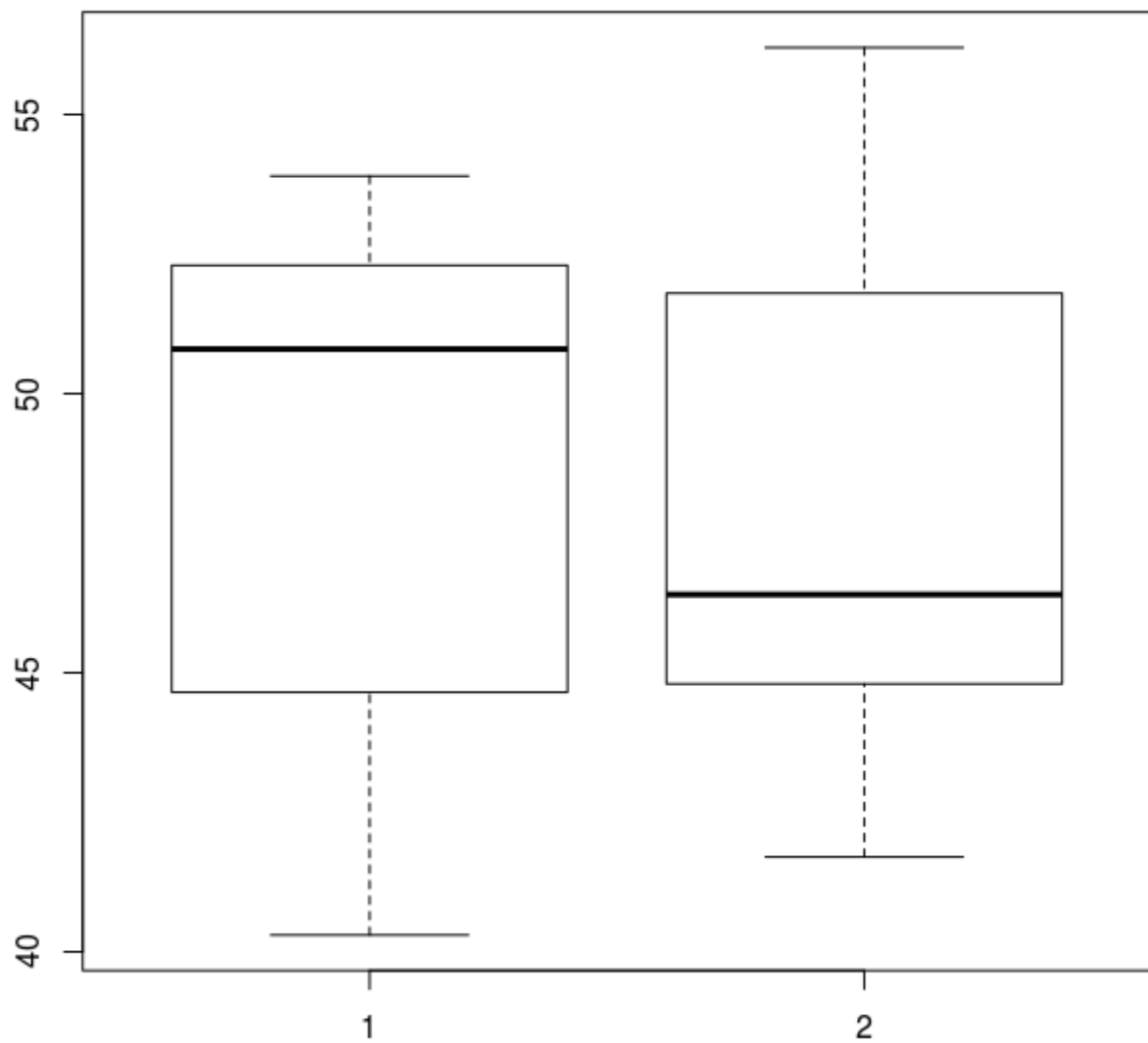


2.3

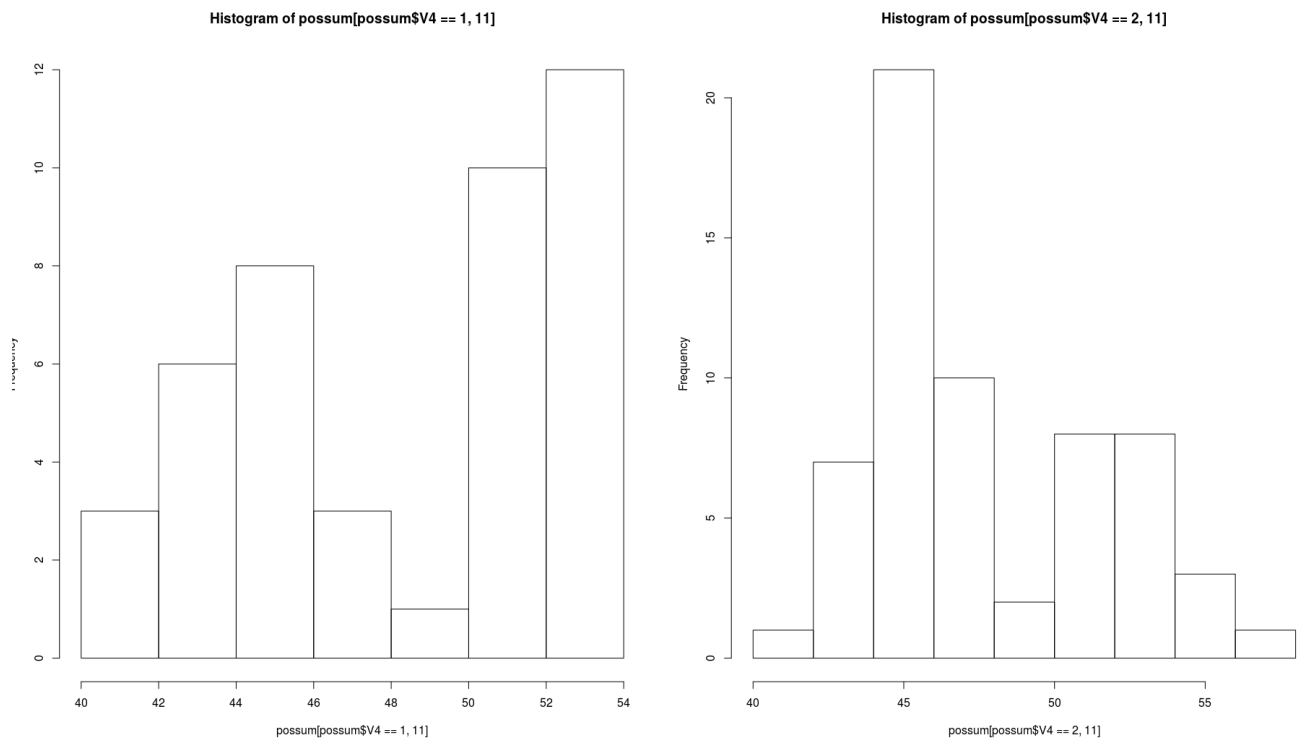
```
> library(DAAG)
> head(possum)
> head(possum)
  case site Pop sex age hdlngth skullw totlngth taill footlngth earconch eye
C3     1    1 Vic  m   8   94.1   60.4    89.0   36.0    74.5    54.5 15.2
C5     2    1 Vic  f   6   92.5   57.6    91.5   36.5    72.5    51.2 16.0
C10    3    1 Vic  f   6   94.0   60.0    95.5   39.0    75.4    51.9 15.5
C15    4    1 Vic  f   6   93.2   57.1    92.0   38.0    76.1    52.2 15.2
C23    5    1 Vic  f   2   91.5   56.3    85.5   36.0    71.0    53.2 15.1
C24    6    1 Vic  f   1   93.1   54.8    90.5   35.5    73.2    53.6 14.2
  chest belly
C3   28.0    36
C5   28.5    33
C10  30.0    34
C15  28.0    34
C23  28.5    33
C24  30.0    32
> with(possum, hist(earconch))
> with(possum, boxplot(earconch~sex))
...
> par(mfrow=c(1,2)) #2 histograms
> hist(possum[possum$sex=="f",11])
> hist(possum[possum$sex=="m",11])
# can also do...
> ?density
> plot(density(possum[possum$sex=="f",11]))
```

Histogram of V11





The measurement distributions have a similar range (with sex 2 having more extreme high values), but with significantly shifted peaks. Sex 1 has its median quite close to the 75th percentile, whereas sex 2 has its peak very close to the 25th percentile. The histogram for 1 would have a tail at the low end, with a peak around 50, and a sharp drop off. The histogram for sex 2 would be similar, but mirrored about $x \approx 47$.



Hmm... I didn't quite guess right. I didn't expect there to be bimodal distributions in the separation by sex based on the box plots. I'm not sure I understand - I would expect that if there were significant numbers of low values for sex 1, then its median would be closer to halfway between the 25th and 75th percentiles. I guess now that I think about it, you wouldn't be able to distinguish the low end distribution accurately just from the box plot.

Note - `density` convolves the data with a kernel, where the bandwidth (width) of the kernel is set in a particular way. Check `?density` for more details.

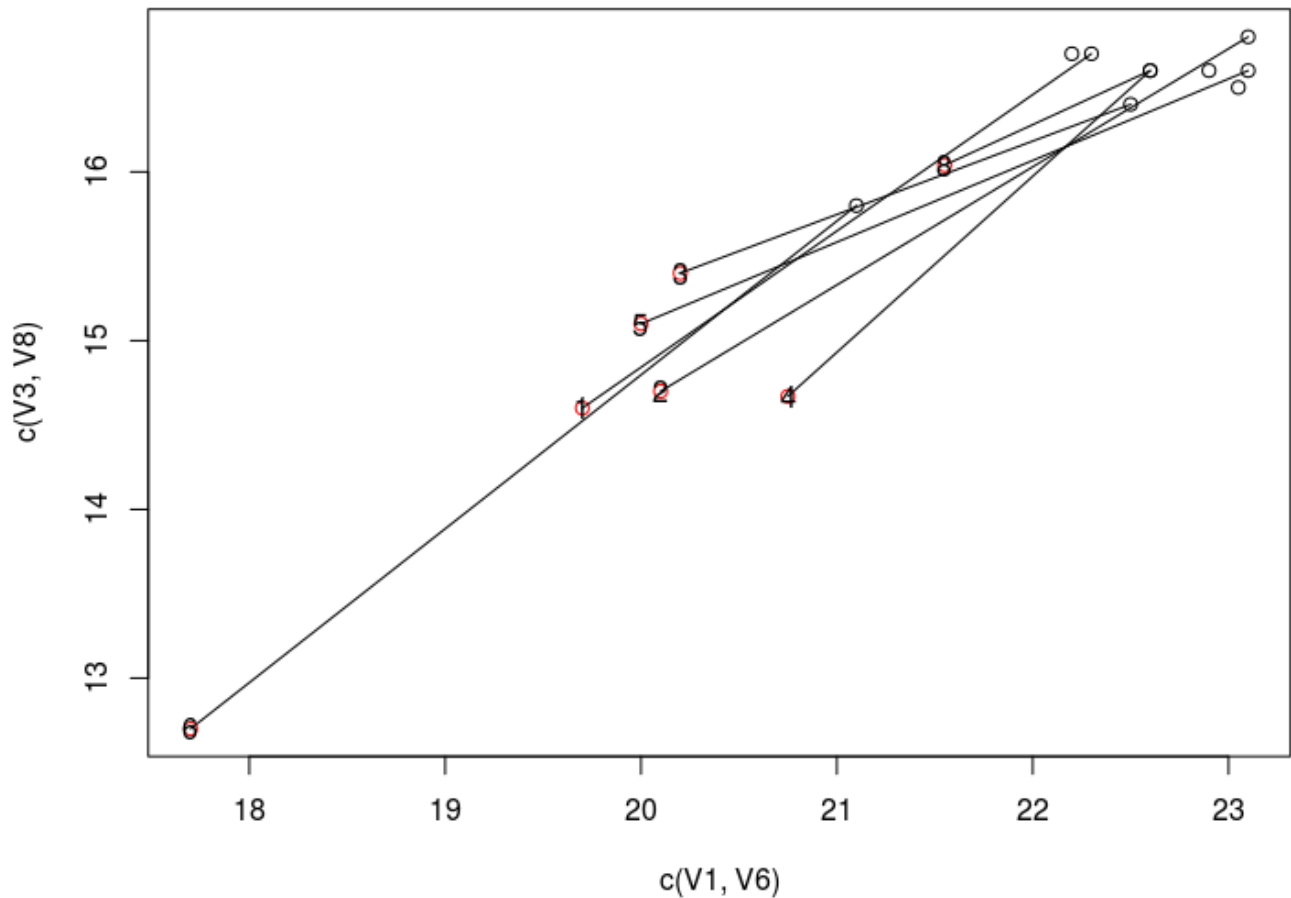
2.5

File format [here](#).

```
> head(cuckoohosts)
      clength cl.sd cbreadth cb.sd cnum hlength hl.sd hbreadth hb.sd
meadow.pipit  22.3  0.89   16.7  0.38  45  19.70  1.25   14.60  0.56
hedge.sparrow  23.1  1.01   16.8  0.52  14  20.10  0.81   14.70 14.70
robin         22.5  0.66   16.4  0.53  16  20.20  0.86   15.40 15.40
wagtails      22.6  0.90   16.6  0.45  26  20.75  1.44   14.67  0.37
tree.pipit    23.1  0.85   16.6  0.44  15  20.00  0.70   15.10  0.48
wren          21.1  0.76   15.8  0.30  15  17.70 17.70   12.70  0.37

      hnum match nomatch
meadow.pipit  74    56     6
hedge.sparrow  26     1    19
robin         57     7    11
wagtails      16    26     3
tree.pipit    27    11     4
wren          NA     0    17
> with(cuckoohosts, plot(c(clength,hlength),c(cbreadth, hbreadth), col=c(rep(1,10),
rep(2,10))))
```

```
> for(i in 1:10)
+   with(cuckoohosts, lines(c(clength[i],hlength[i]),c(cbreadth[i], hbreadth[i])))
> with(cuckoohosts, text(hlength, hbreadth, rownames(cuckoohosts)))
```



A long line implies that the cuckoo eggs in that particular nest were very different (length and/or breadth-wise) compared to the host eggs. Short lines indicate cuckoo eggs that are quite similar (on average) to the host bird's eggs.

2.10

I believe this dataset is described [here](#).

```
> library(MASS)
> head(Animals)

      body brain
Mountain beaver    1.35   8.1
Cow              465.00 423.0
Grey wolf         36.33 119.5
Goat              27.66 115.0
Guinea pig        1.04   5.5
Dipliodocus     11700.00  50.0
> with(Animals, cor(brain, body))
```

```
[1] -0.005341163
> with(animals, cor(log(brain), log(body)))
[1] 0.7794935
> with(animals, cor(log(brain), log(body), method="spearman"))
[1] 0.7162994
> par(mfrow=c(1,2))
> with(animals, plot(body,brain))
> with(animals, plot(log(body),log(brain)))
```

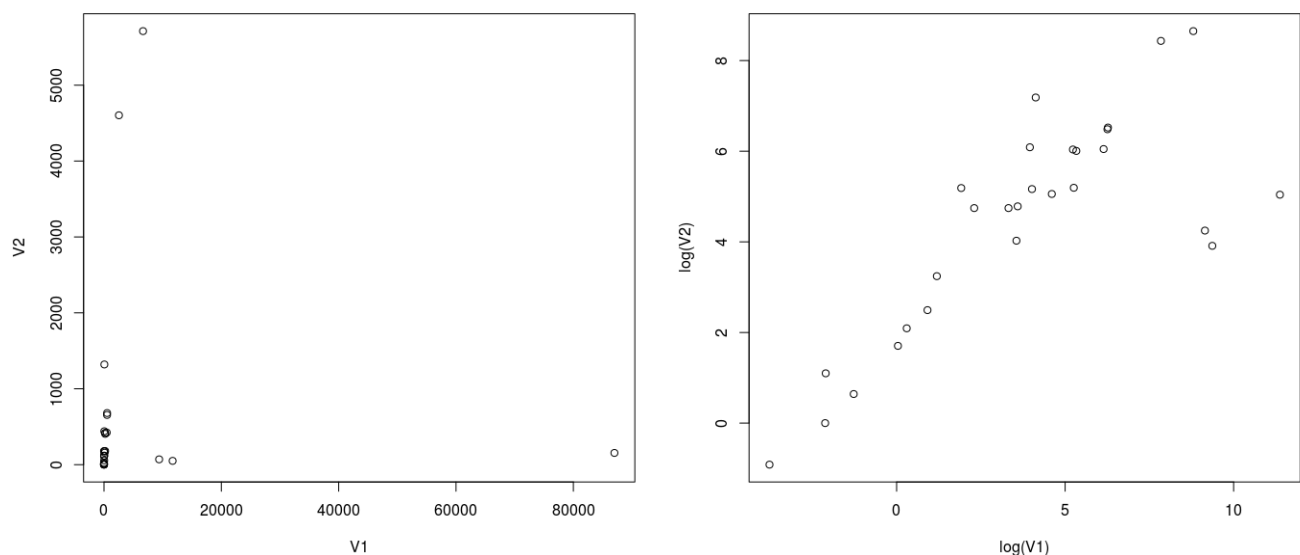
From the documentation:

'var', 'cov' and 'cor' compute the variance of 'x' and the covariance or correlation of 'x' and 'y' if these are vectors. If 'x' and 'y' are matrices then the covariances (or correlations) between the columns of 'x' and the columns of 'y' are computed.

'cov2cor' scales a covariance matrix into the corresponding correlation matrix *efficiently*.

...

For 'cor()', if 'method' is "'kendall'" or "'spearman'", Kendall's tau or Spearman's rho statistic is used to estimate a rank-based measure of association. These are more robust and have been recommended if the data do not necessarily come from a bivariate normal distribution. For 'cov()', a non-Pearson method is unusual but available for the sake of completeness. Note that "'spearman'" basically computes 'cor(R(x), R(y))' (or 'cov(., .)') where 'R(u) := rank(u, na.last = "keep")'. In the case of missing values, the ranks are calculated depending on the value of 'use', either based on complete observations, or based on pairwise completeness with reranking for each pair.



A [Bivariate Normal](#) distribution is just something that is normal in 2 dimensions. This data certainly doesn't look normal to me, so I would guess that the Spearman correlation of the log data is the most appropriate here.

Note - logging the data shouldn't change the spearman correlation.

```
> with(Animals, cor(brain, body, method="spearman"))
[1] 0.7162994
> with(Animals, cor(log(brain), log(body), method="spearman"))
[1] 0.7162994
```

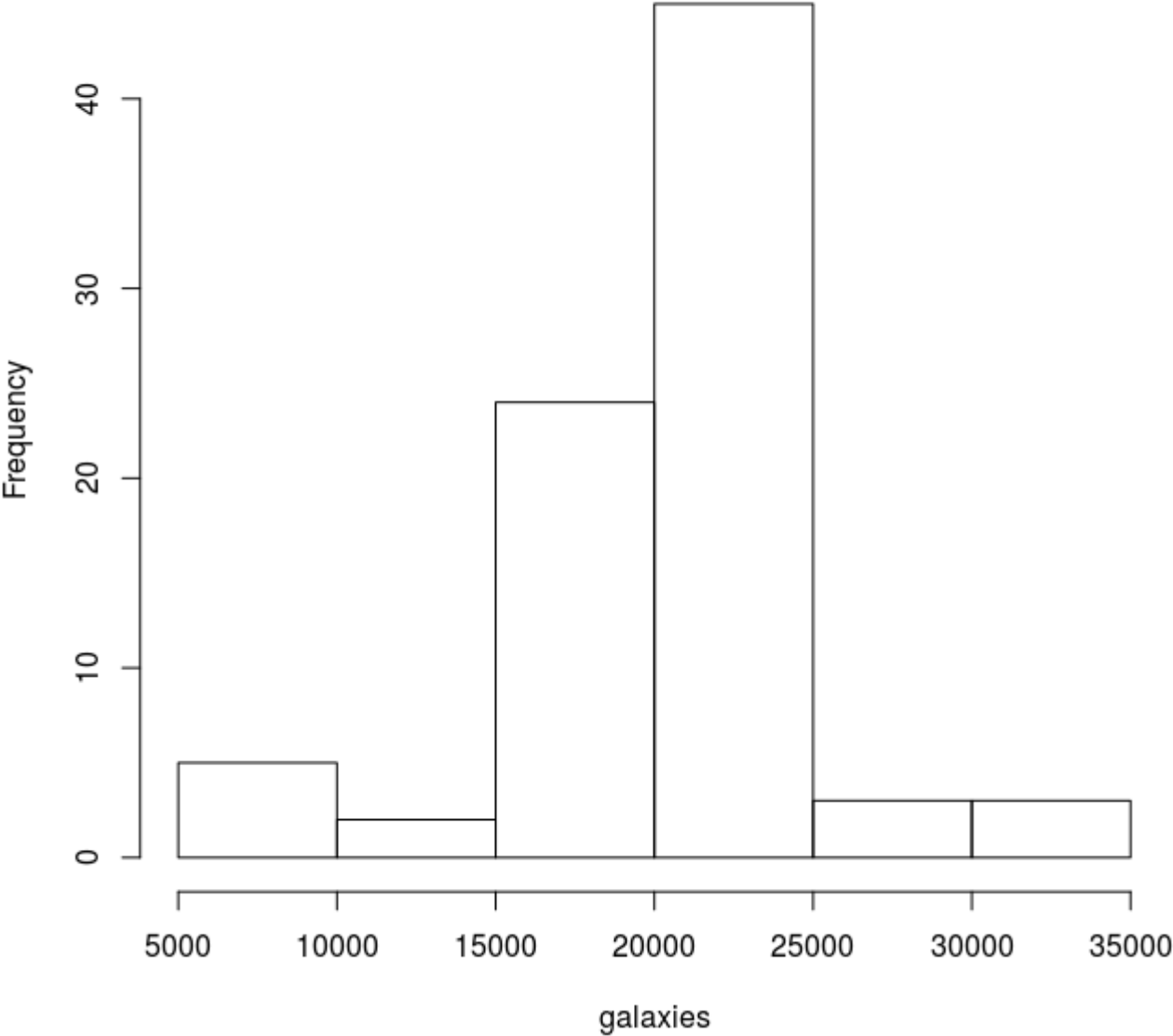
Another way to look at this - as long as the mapping from $x,y \rightarrow x',y'$ is monotonic, the spearman correlation will be the same before and after the correlation.

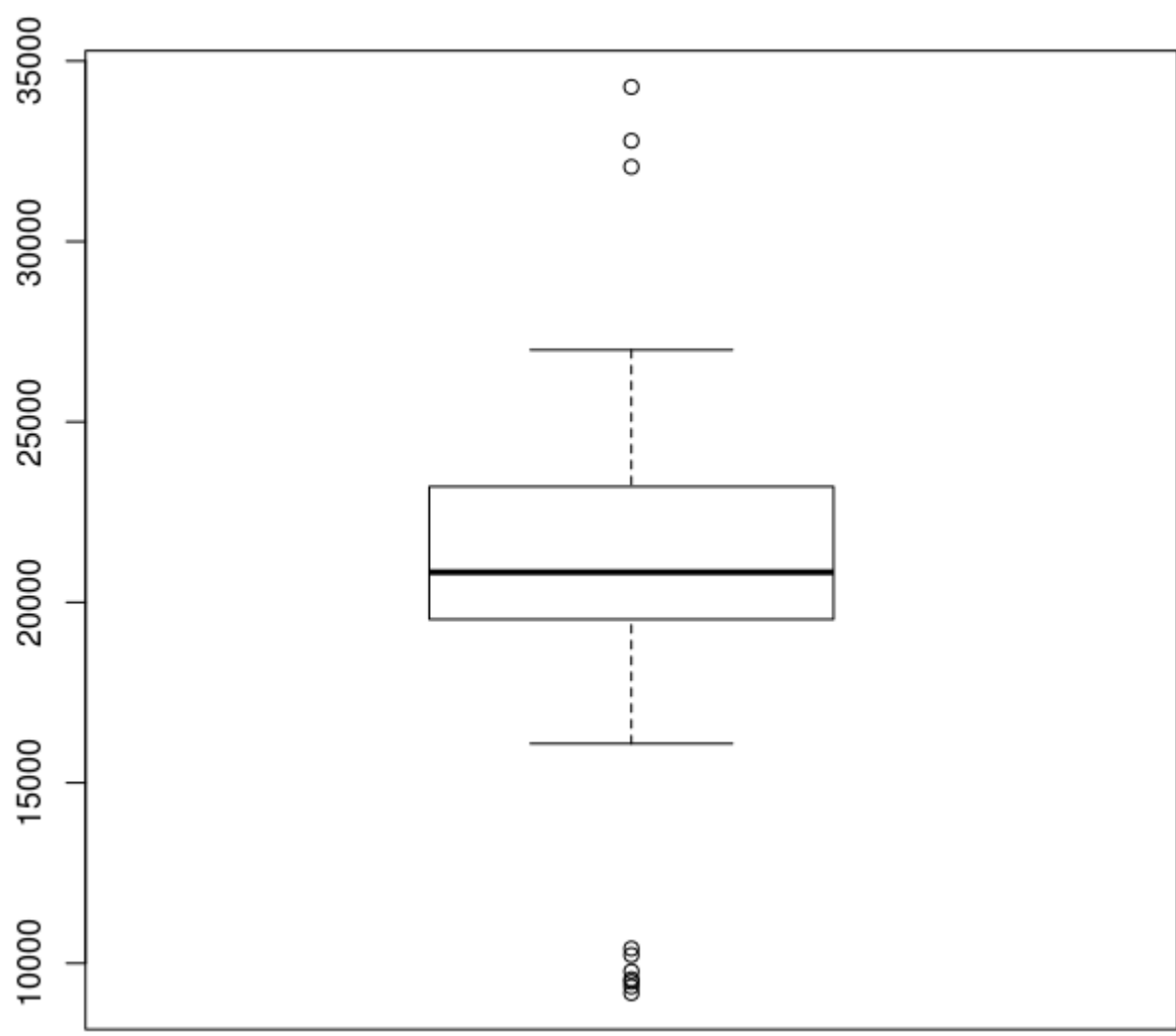
Rank correlation is just correlating the underlying data distribution. Note that a pearson correlation is just the covariance / $(\text{var}(x1) * \text{var}(x2))$, which is measuring the extent to which a bivariate distribution is oblong (e.g. a circle is not correlated, but if you flatten it to a line, it's perfectly correlated).

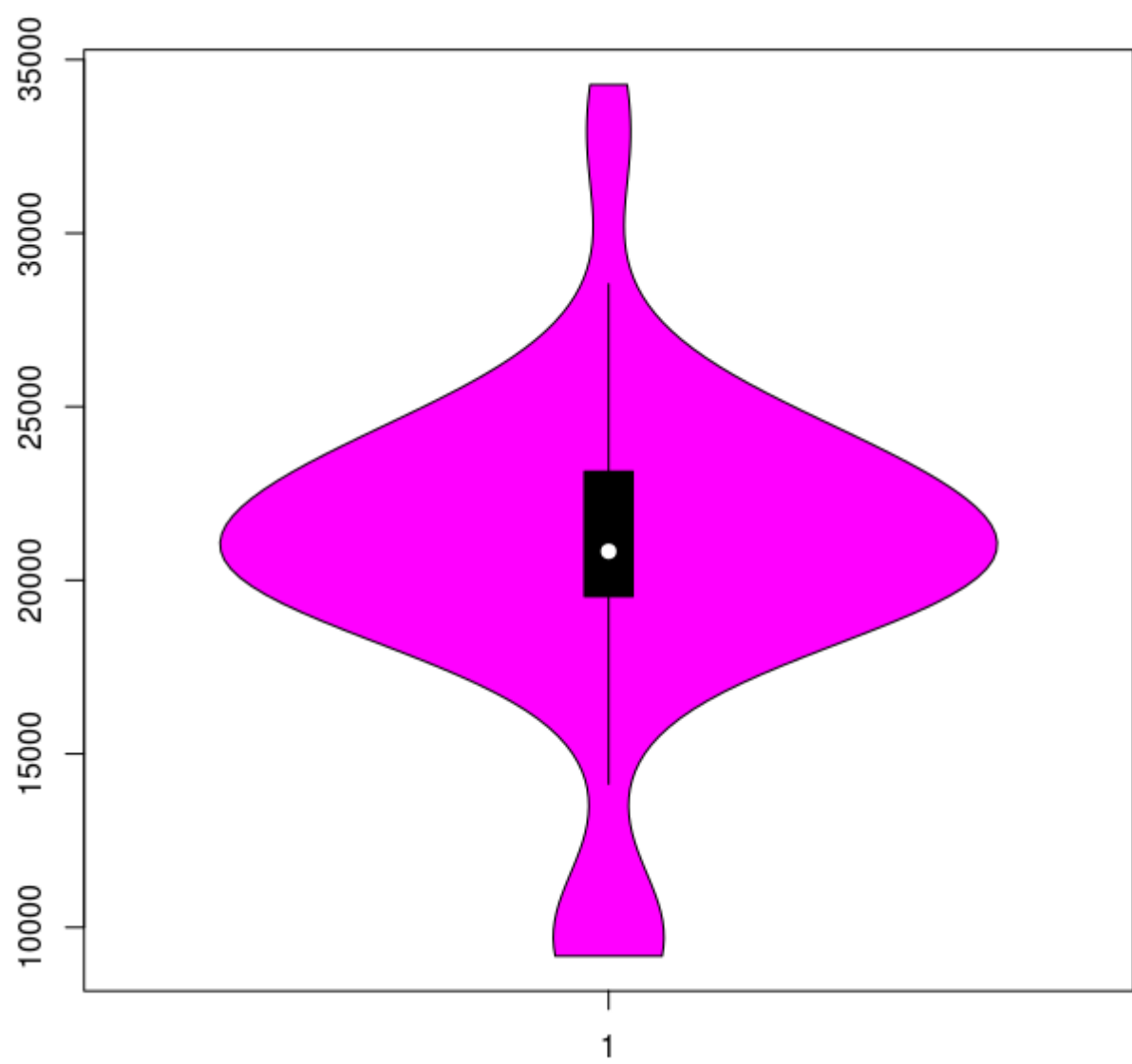
2.13

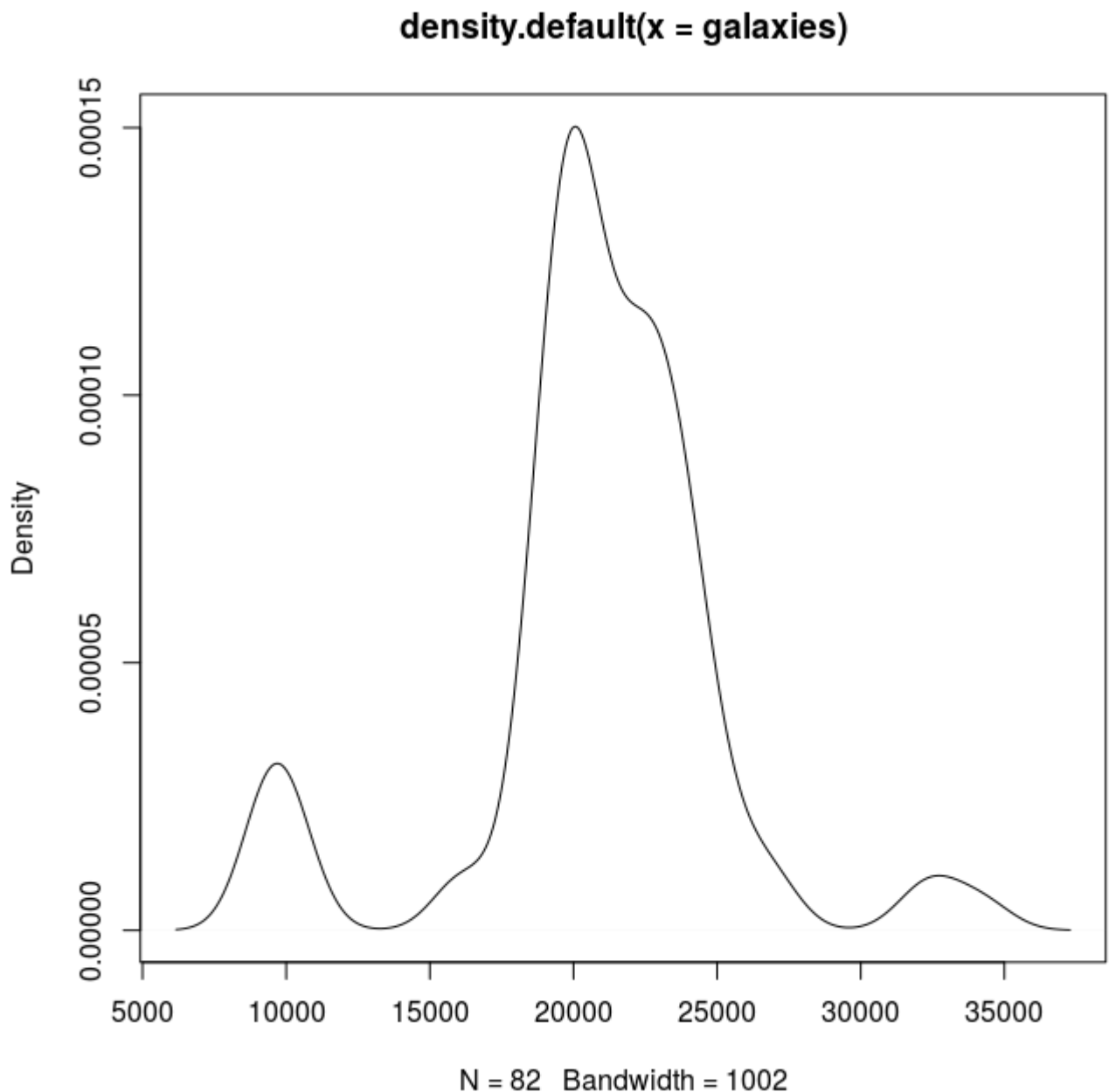
```
> library(MASS)
> library(vioplot)
> head(galaxies)
[1] 9172 9350 9483 9558 9775 10227
> hist(galaxies)
> boxplot(galaxies)
> vioplot(galaxies)
> plot(density(galaxies))
# see ?galaxies for another example of plotting the data
```

Histogram of galaxies









The distribution does not look skewed to me. It seems quite clustered around ~22,000. There is something interesting going on around 10,000 - a lot of low end "outliers". This could perhaps be a 2nd population that hasn't been sampled well enough to see smaller distribution below 10,000 (and in fact could be observation bias if it's more difficult to see lower speed galaxies).

Preliminary Exploratory Data Analysis
