

# Lab 10 Solutions

1) Calculate the mean and variance of each column, by using `apply()` function.

```
states <- row.names(USArrests)
states

## [1] "Alabama"      "Alaska"      "Arizona"     "Arkansas"
## [5] "California"   "Colorado"    "Connecticut" "Delaware"
## [9] "Florida"      "Georgia"     "Hawaii"      "Idaho"
## [13] "Illinois"     "Indiana"     "Iowa"        "Kansas"
## [17] "Kentucky"    "Louisiana"   "Maine"       "Maryland"
## [21] "Massachusetts" "Michigan"    "Minnesota"   "Mississippi"
## [25] "Missouri"    "Montana"     "Nebraska"    "Nevada"
## [29] "New Hampshire" "New Jersey" "New Mexico"  "New York"
## [33] "North Carolina" "North Dakota" "Ohio"        "Oklahoma"
## [37] "Oregon"      "Pennsylvania" "Rhode Island" "South Carolina"
## [41] "South Dakota" "Tennessee"   "Texas"       "Utah"
## [45] "Vermont"     "Virginia"    "Washington"  "West Virginia"
## [49] "Wisconsin"   "Wyoming"

names(USArrests)

## [1] "Murder" "Assault" "UrbanPop" "Rape"

apply(USArrests, 2, mean)

## Murder Assault UrbanPop Rape
## 7.788 170.760 65.540 21.232

apply(USArrests, 2, var)

## Murder Assault UrbanPop Rape
## 18.97047 6945.16571 209.51878 87.72916
```

2) What conclusions can you draw from 1)? And consequently what transformation would you do to your dataset?

We see that there are on average three times as many rapes as murders, and more than eight times as many assaults as rapes.

Not surprisingly, the variables also have vastly different variances: the UrbanPop variable measures the percentage of the population in each state living in an urban area, which is not a comparable number to the number of rapes in each state per 100,000 individuals.

If we failed to scale the variables before performing PCA, then most of the principal components that we observed would be driven by the Assault variable, since it has by far the largest mean and variance. Thus, it is important to standardize the variables to have mean zero and standard deviation one before performing PCA.

3) Perform principal component analysis using the `prcomp()` function.

```
pr.out <- prcomp(USArrests, scale = TRUE)
```

4) Check the results, report the number of PCs and their center, scale, and rotation.

```
names(pr.out)
```

```
## [1] "sdev"      "rotation" "center"    "scale"     "x"
```

```
pr.out$center  #Don't write it as "centre"!
```

```
## Murder Assault UrbanPop Rape
## 7.788 170.760 65.540 21.232
```

```
pr.out$scale
```

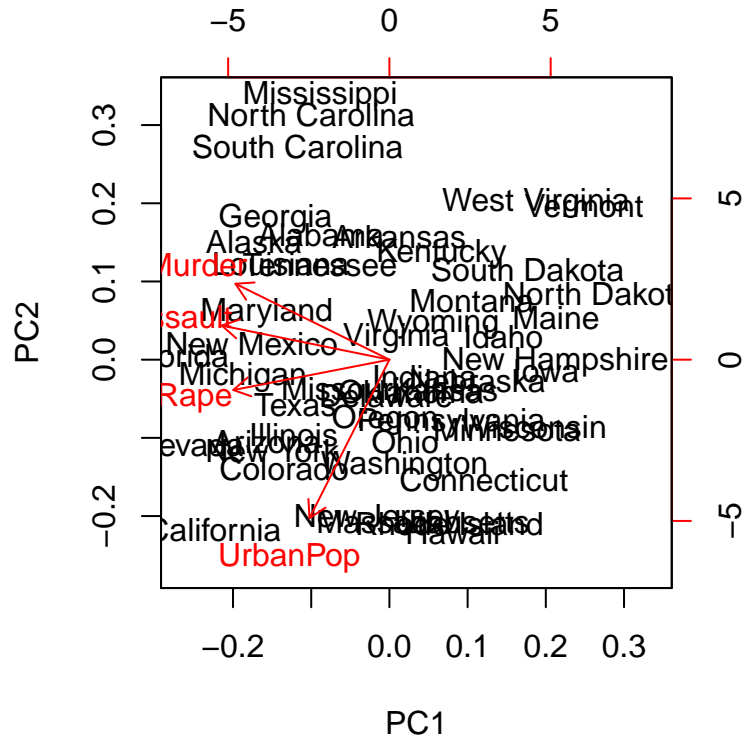
```
## Murder Assault UrbanPop Rape
## 4.355510 83.337661 14.474763 9.366385
```

```
pr.out$rotation
```

```
##          PC1          PC2          PC3          PC4
## Murder -0.5358995  0.4181809 -0.3412327  0.64922780
## Assault -0.5831836  0.1879856 -0.2681484 -0.74340748
## UrbanPop -0.2781909 -0.8728062 -0.3780158  0.13387773
## Rape     -0.5434321 -0.1673186  0.8177779  0.08902432
```

5) Plot the first two PCs.

```
biplot(pr.out, scale = TRUE)
```

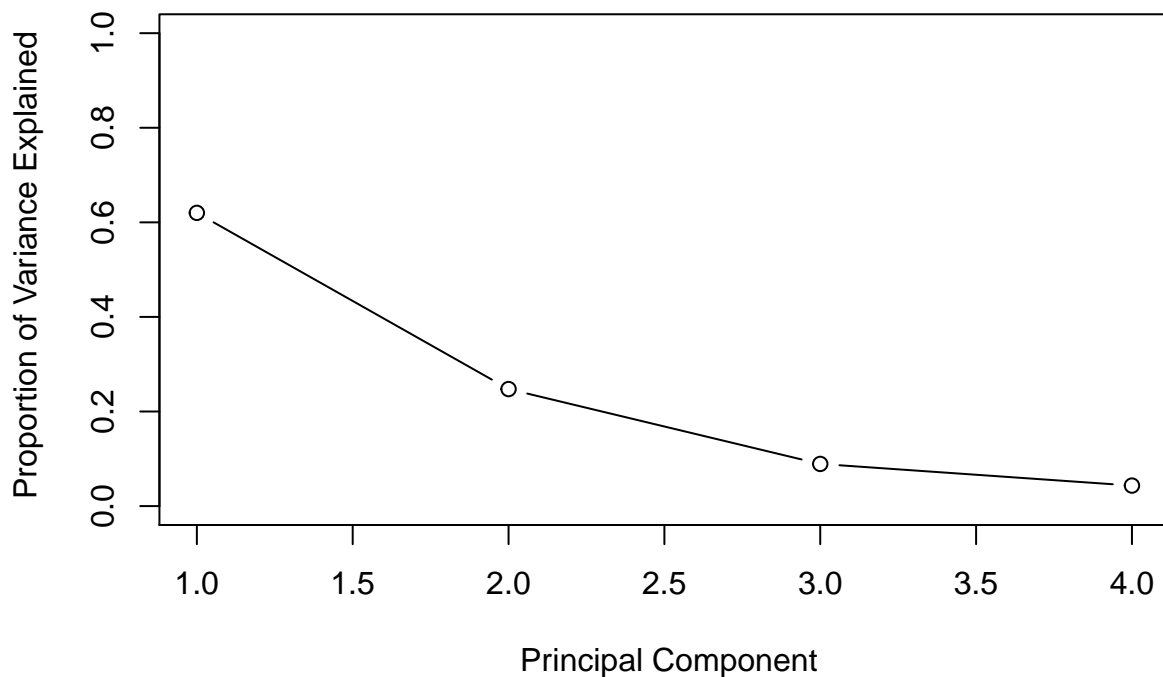


6) What are the standard deviation of each principal component? Based on this result, calculate the variance explained by each PC and the proportion of variance explained by each PC.

```
pr.out$sdev
## [1] 1.5748783 0.9948694 0.5971291 0.4164494
pr.var <- pr.out$sdev^2
pr.var
## [1] 2.4802416 0.9897652 0.3565632 0.1734301
pve <- pr.var/sum(pr.var)
pve
## [1] 0.62006039 0.24744129 0.08914080 0.04335752
```

7) Plot the PVE explained by each component as well as the cumulative PVE.

```
plot(pve,
     xlab = "Principal Component", ylab = "Proportion of Variance Explained",
     ylim=c(0,1), type='b')
```



```
plot(cumsum(pve),
     xlab = "Principal Component ", ylab = "Cumulative Proportion of Variance Explained",
     ylim = c(0,1), type = 'b')
```

