# Principle component analysis

Evan Cummings

CSCI 548 – Douglas W. Raiford – Pattern Recognition

November 16, 2016

## 1 `class` **command :**

The class of the `iris` dataset is a `data.frame`, a "tightly coupled collections of variables which share many of the properties of matrices and of lists, used as the fundamental data structure by most of R's modeling software."



## 2 `summary` **command :**

The `summary` command applied to the `iris` dataset, is a "generic function used to produce result summaries of the results of various model fitting functions. The function invokes particular methods which depend on the 'class' of the first argument."



## 3 `labels` **command :**

The `labels` command: "Find[s] a suitable set of labels from an object for use in printing or plotting. For example, a generic function."



## 4 `colnames` **command :**

The `colnames` command: "Retrieve[s] or set[s] the row or column names of a matrix-like object."



## 5 **tab-completion :**

Pressing `tab` after a `data.frame` object with accessor-operator `$` will provide variable names.

# 6   Iris data PCA plotting :







Figure 1: PCA-plot showing what appears to be two clusters, despite the fact that there are three distinct species of iris.

The `plot` command, a "Generic function for plotting of R objects. For more details about the graphical parameter arguments, see 'par'.

For simple scatter plots, `plot.default` will be used. However, there are `plot` methods for many R objects, including `function`'s, `data.frame`'s, `density` objects, etc. Use `methods(plot)` and the documentation for these."

**Proportions of Variance**



**Fruit Data**



## 6.1 Source code :

```
# store the iris 'classes' :
c = iris[,5]

# store the iris 'data' :
d = iris[,seq(1,4)]

# perform PCA on d :
p = prcomp(d)

# set the 1st principle component :
x = p$x[,1]

# set the 2nd principle component :
y = p$x[,2]

# plot them :
png('../doc/images/pca_plot.png')
plot(x,y)
dev.off()

# plot them in red :
png('../doc/images/pca_plot_red.png')
plot(x,y, col='red', bg='red', type='p', pch=21)
dev.off()

# store indexes of classes :
s  = which(c == 'setosa')
vc = which(c == 'versicolor')
vg = which(c == 'virginica')

# color vector :
cc     = as.vector(c)
cc[s]  = 'red'
cc[vc] = 'blue'
cc[vg] = 'green'

# plot them colored by class :
png('../doc/images/pca_plot_class.png')
X11(width=4, height=4)
par(mar=c(2.5,2.5,2.5,.1),mgp = c(1.5, .5, 0))
plot(x,y, col=cc, bg=cc, type='p', pch=21, xlab='First Principal Component',
     ylab='Second Principal Component', main='Iris Data')
legend('topright', levels(c),
       col    = c('red', 'green', 'blue'),
       pt.bg = c('red', 'green', 'blue'), pch = 21)
dev.off()

# do a barplot of the proportions of variance :
sig  = p$sdev**2
pvar =  sig / sum(sig)
png('../doc/images/pca_pvar_barplot.png')
barplot(pvar, names.arg=c('PC1','PC2','PC3','PC4'),
        main='Proportions of Variance')
dev.off()
```

# 7 Fruit data :

While the iris data appears to be highly grouped, the fruit data appears much less so. However, the lemons seem to be clustered distintly from the other fruit. Therefore, I believe that these data *do* show enough structure for machine learning techniques; with apples, oranges, and peaches potentially presenting the highest challenge.
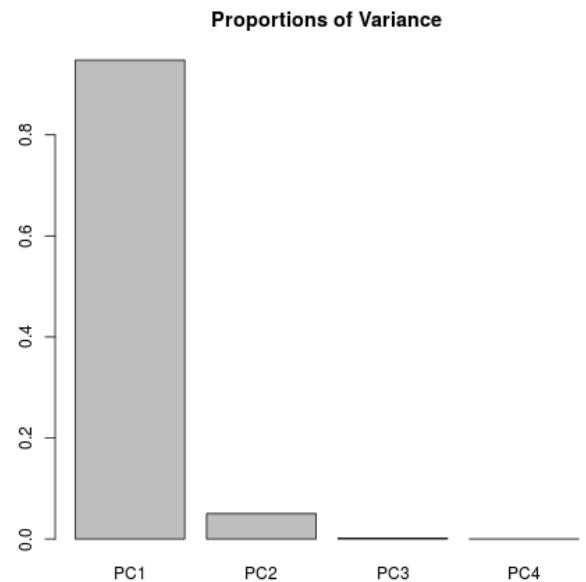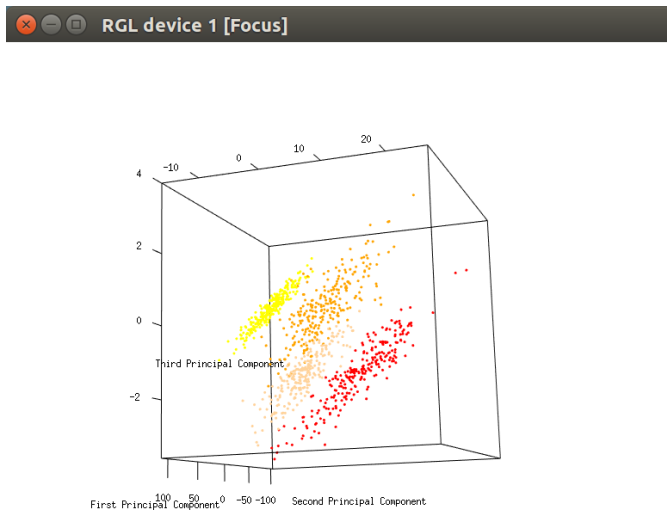
**Proportions of Variance**



Figure 2: The proportions of the variance in the fruit data. The principle component captures ≈ 94% of the variation.

```
o = which(c == 'orange')
z = which(c == 'peach')

# color vector :
cols  = c('red', 'yellow', 'orange', 'burlywood1')
cc    = as.vector(c)
cc[a] = cols[1]
cc[l] = cols[2]
cc[o] = cols[3]
cc[z] = cols[4]

# plot them colored by class :
png('../doc/images/pca_plot_fruit.png')
#X11(width=4, height=4)
par(mar=c(2.5,2.5,2.5,.1),mgp = c(1.5, .5, 0))
plot(x,y, col=cc, bg=cc, type='p', pch=21, xlab='First Principal Component',
     ylab='Second Principal Component', main='Fruit Data')
legend('topright', levels(c), col = cols, pt.bg = cols, pch = 21)
dev.off()

# do a barplot of the proportions of variance :
sig  = p$sdev**2
pvar =  sig / sum(sig)
png('../doc/images/pca_pvar_fruit_barplot.png')
barplot(pvar, names.arg=c('PC1','PC2','PC3','PC4'),
        main='Proportions of Variance')
dev.off()

# import the rgl library for 3D stuff:
library(rgl)

plot3d(x,y,z, col=cc, bg=cc, type='p', pch=21,
       xlab='First Principal Component',
       ylab='Second Principal Component',
       zlab='Third Principal Component')

# plot them colored by class :
png('../doc/images/pca_plot_fruit_23.png')
#X11(width=4, height=4)
par(mar=c(2.5,2.5,2.5,.1),mgp = c(1.5, .5, 0))
plot(y,z, col=cc, bg=cc, type='p', pch=21, xlab='Second Principal Component',
     ylab='Third Principal Component', main='Fruit Data')
legend('topright', levels(c), col = cols, pt.bg = cols, pch = 21)
dev.off()
```

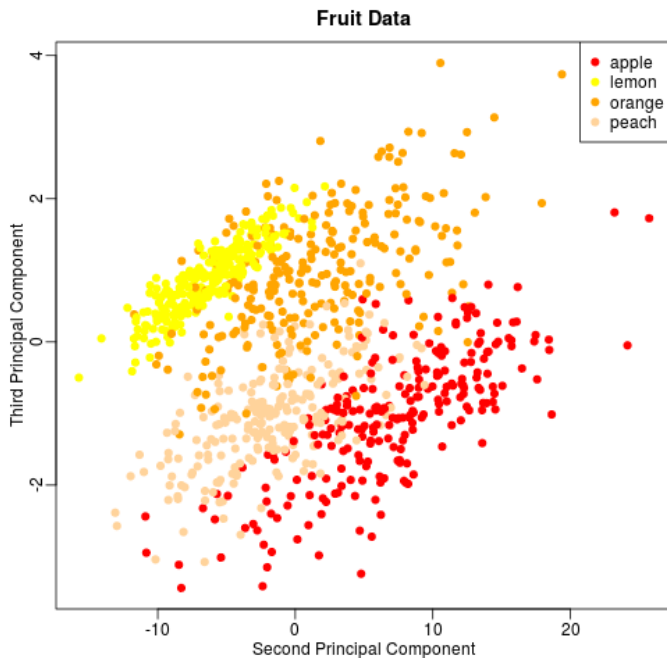Figure 3: In 3D, we can see much more distinct clustering.



Figure 4: The second and third principle components, highlighting the clustering nature of the fruit data.

## 7.1 Source code :

```
# read the variable back in :
f = read.csv("../data/fruit.csv")

# store the fruit 'classes' :
c = f[,5]

# store the fruit 'data' :
d = f[,seq(1,4)]

# perform PCA on d :
p = prcomp(d)

# set the 1st principle component :
x = p$x[,1]

# set the 2nd principle component :
y = p$x[,2]

# set the 3nd principle component :
z = p$x[,3]

# store indexes of classes :
a = which(c == 'apple')
l = which(c == 'lemon')
```