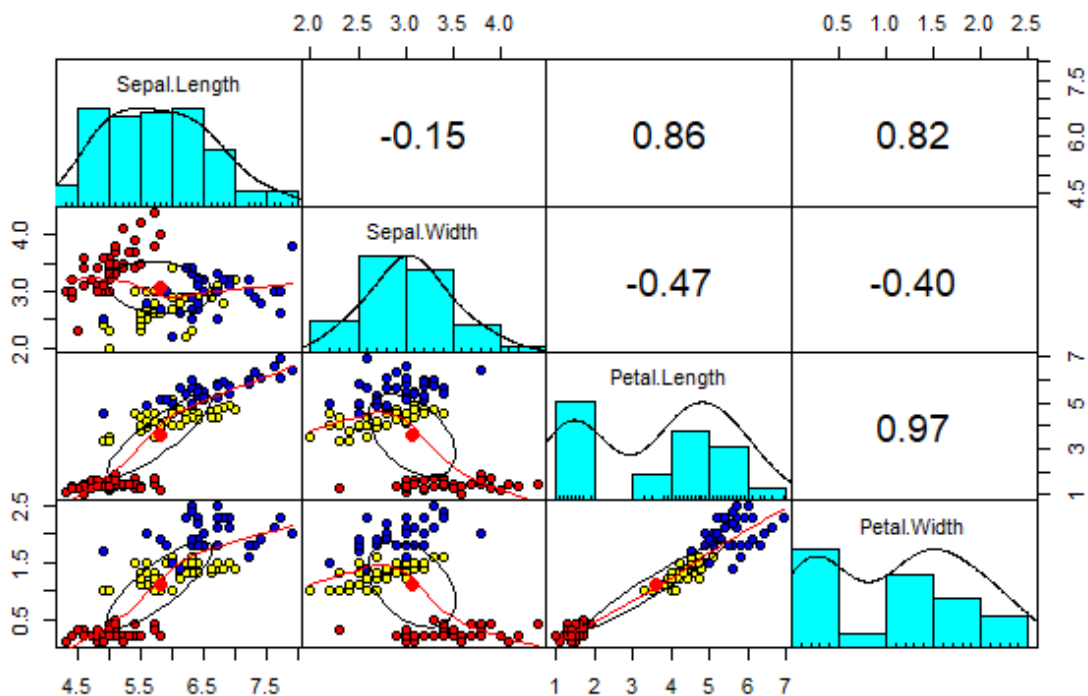PCA，主成分分析是一种有效的降维方法，常用于高维数据的处理，如基因表达谱数据。本教程主要展示PCA在R语言中的实现和可视化。

# 数据准备

所用数据为**iris**数据集。

```r
# data
data(iris)
str(iris)

# partition data # 将数据划分成训练集和测试集，一把的分析不用划分
set.seed(111)
ind = sample(2, nrow(iris),replace = T, prob = c(.8,.2))
train = iris[ind == 1,]
test = iris[ind == 2,]

# scatter plot and correlations 绘图展示数据集变量之间的关系
library(psych)
pairs.panels(train[,-5],
             gap = 0,
             bg = c('red','yellow','blue')[train$Species],
             pch = 21)
```



# PCA计算

```r
# PCA
```

```
2   pc = prcomp(train[,-5],# 数据最后一列是分类变量，不选择
3               center = T,# 数据中心化
4               scale. = T) #数据标准化
5   pc$center
6   pc$scale
7   print(pc)
8
9   > print(pc)
10  Standard deviations (1, .., p=4):
11  [1] 1.7173318 0.9403519 0.3843232 0.1371332
12
13  Rotation (n x k) = (4 x 4):
14                     PC1         PC2         PC3         PC4
15  Sepal.Length  0.5147163 -0.39817685  0.7242679  0.2279438
16  Sepal.Width  -0.2926048 -0.91328503 -0.2557463 -0.1220110
17  Petal.Length  0.5772530 -0.02932037 -0.1755427 -0.7969342
18  Petal.Width   0.5623421 -0.08065952 -0.6158040  0.5459403
```

查看**PC**的解释百分百：

```
1   > summary(pc)
2   Importance of components:
3                            PC1     PC2     PC3     PC4
4   Standard deviation     1.7173  0.9404  0.38432  0.1371
5   Proportion of Variance 0.7373  0.2211  0.03693  0.0047
6   Cumulative Proportion  0.7373  0.9584  0.99530  1.0000
```
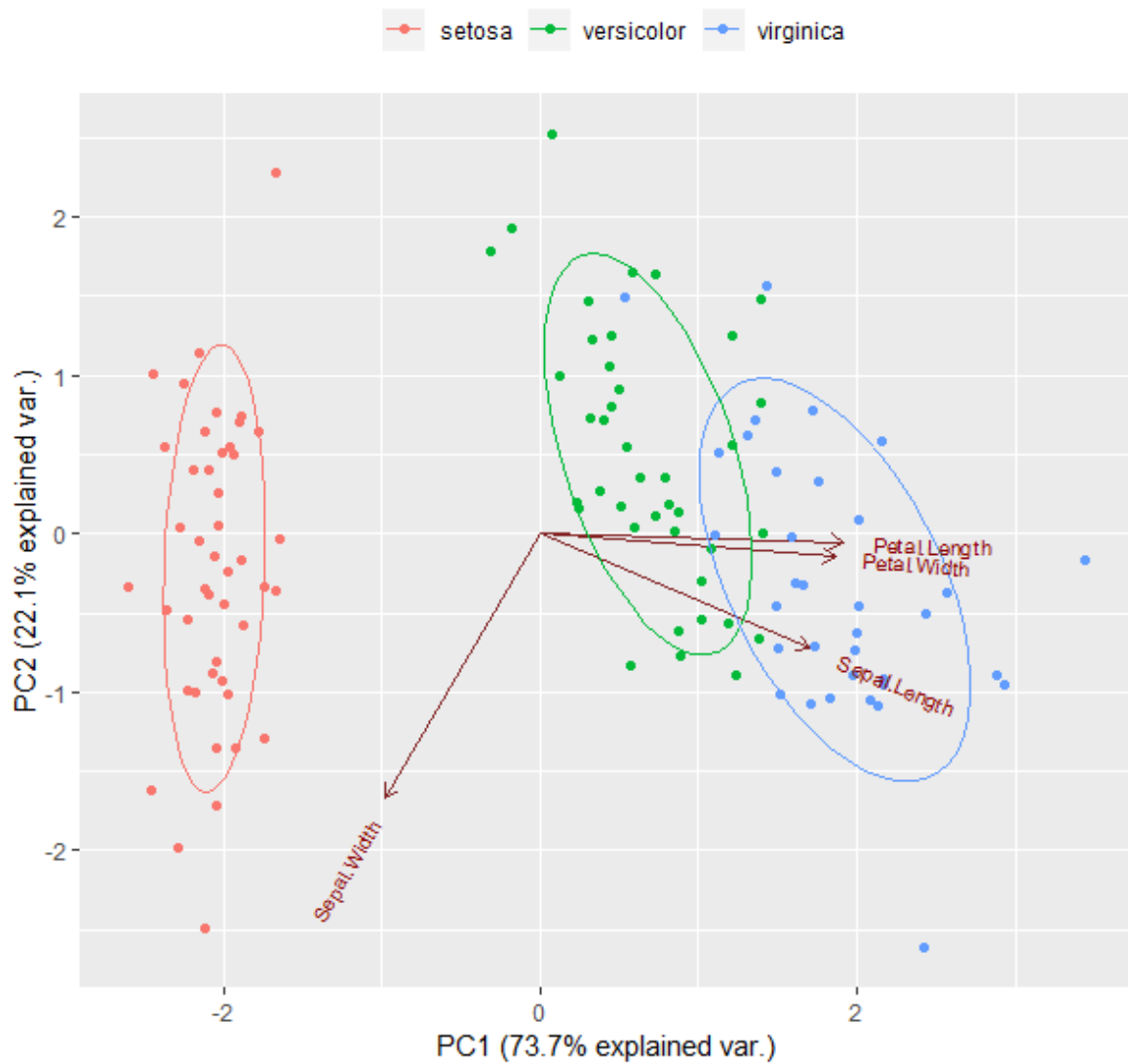
直接看**Proportion of Variance**即可。

# 可视化

```
1   # bi-plot
2   library(ggbiplot)
3   g = ggbiplot(pc,
4               obs.scale = 1,
5               var.scale = 1,
6               groups = train$Species,
7               ellipse = T, # 置信区间椭圆
8               circle = F,
9               ellipse.prob = 0.68) # 椭圆覆盖多少数据点
10  g = g + scale_color_discrete(name = '') + theme(legend.direction =
    'horizontal',
11                                                 legend.position = 'top')
12  print(g)
```

# 模型预测

一般是不会涉及到模型预测的。

```r
# prediction with Principal Components
trg = predict(pc, train)
trg = data.frame(trg, train$Species)

tst = predict(pc,test)
tst = data.frame(tst, test$Species)

# multinomial logistic regression with first two PCs
library(nnet)
trg$species = relevel(trg$train.Species, ref = 'setosa')
mymodel = multinom(train.Species~PC1+PC2,data = trg)
summary(mymodel)

# Confusion matrix
p = predict(mymodel, tst)
tab = table(p, test$Species)
tab
```

```
18   1-sum(diag(tab))/sum(tab)
```

```
1   > tab
2
3   p              setosa versicolor virginica
4     setosa            5          0         0
5     versicolor        0          9         3
6     virginica         0          1        12
7   > 1-sum(diag(tab))/sum(tab)
8   [1] 0.1333333
```

**致谢：** 感谢YouTube博主*Bharatendra Rai*博士的视频！