```r
library(cluster)
library(factoextra)
library(NbClust)
library(fpc)
library(ggplot2)
library(ggfortify)

data(iris)
head(iris)
iris2 <- iris[1:4]
```

# PCA

```r
> PCAtemp <- prcomp(iris2)
> summary(PCAtemp)
Importance of components:
                          PC1     PC2    PC3     PC4
Standard deviation     2.0563 0.49262 0.2797 0.15439
Proportion of Variance 0.9246 0.05307 0.0171 0.00521
Cumulative Proportion  0.9246 0.97769 0.9948 1.00000
> PCAtemp$rotation
                     PC1         PC2         PC3        PC4
Sepal.Length  0.36138659 -0.65658877  0.58202985  0.3154872
Sepal.Width  -0.08452251 -0.73016143 -0.59791083 -0.3197231
Petal.Length  0.85667061  0.17337266 -0.07623608 -0.4798390
Petal.Width   0.35828920  0.07548102 -0.54583143  0.7536574
```
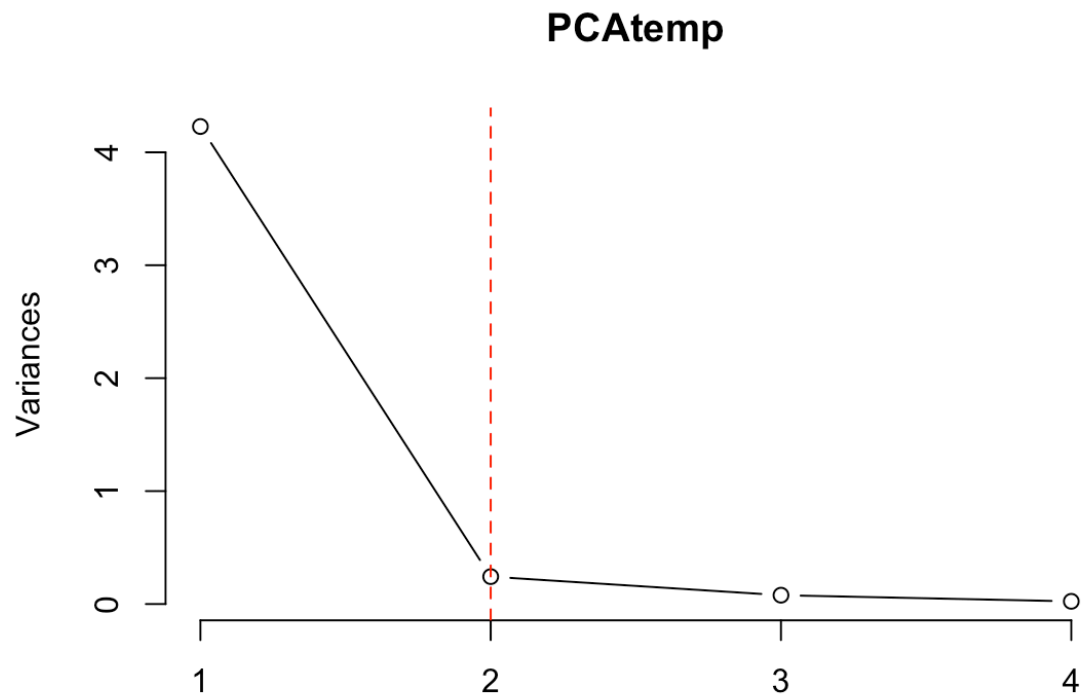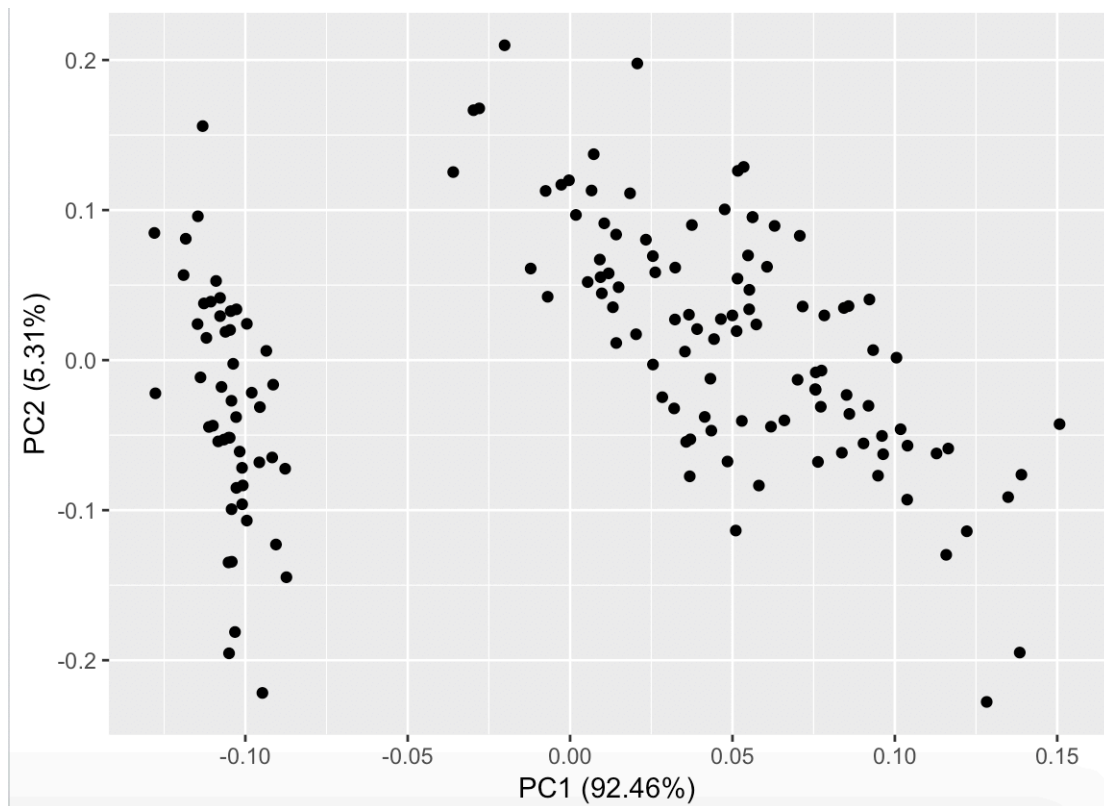
The summary shows that the first two Principal Components account for most of the variation in all four components. "Petal. Length" has a large weight in principal component 1, while "Sepal. Width" has a small weight. The second principal component is mostly "Sepal. Width" and "Sepal. Length" with small weightings on the other variables.
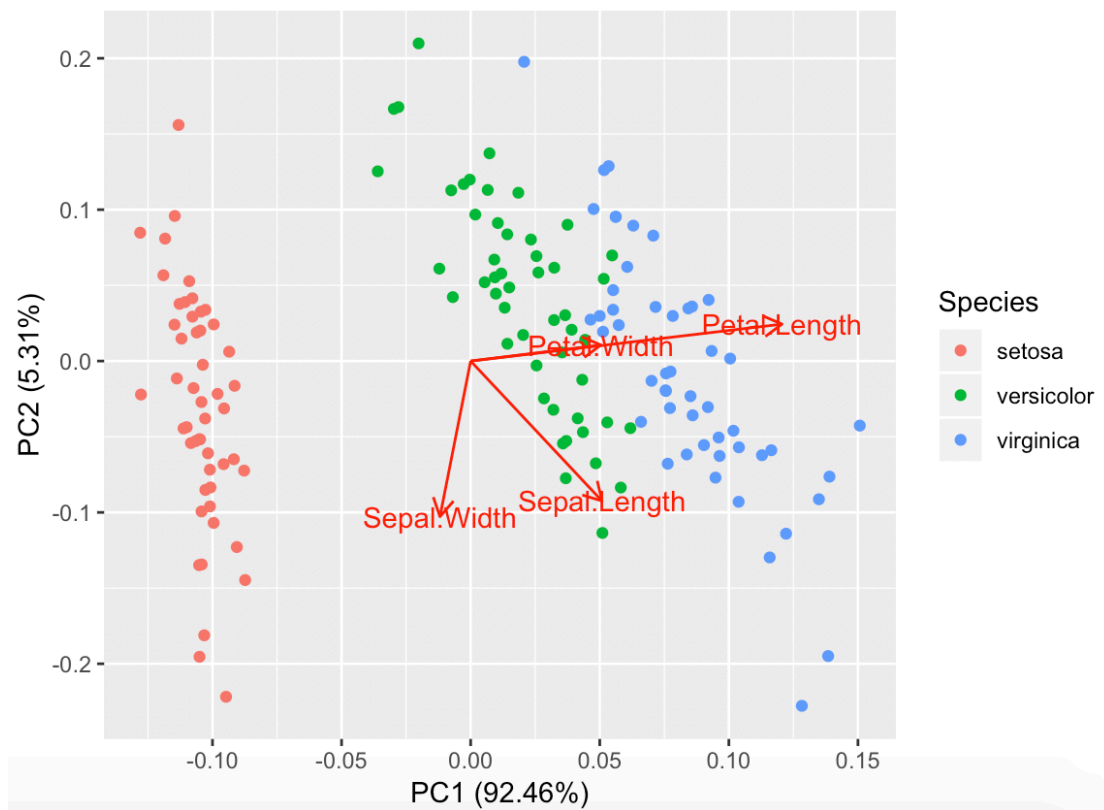
```r
screeplot(PCAtemp, type = "lines")
abline(v = 2,lty =2, col = "red")

autoplot(PCAtemp)
autoplot(PCAtemp, data = iris, col = "Species", loadings = TRUE, loadings.label = TRUE)
```
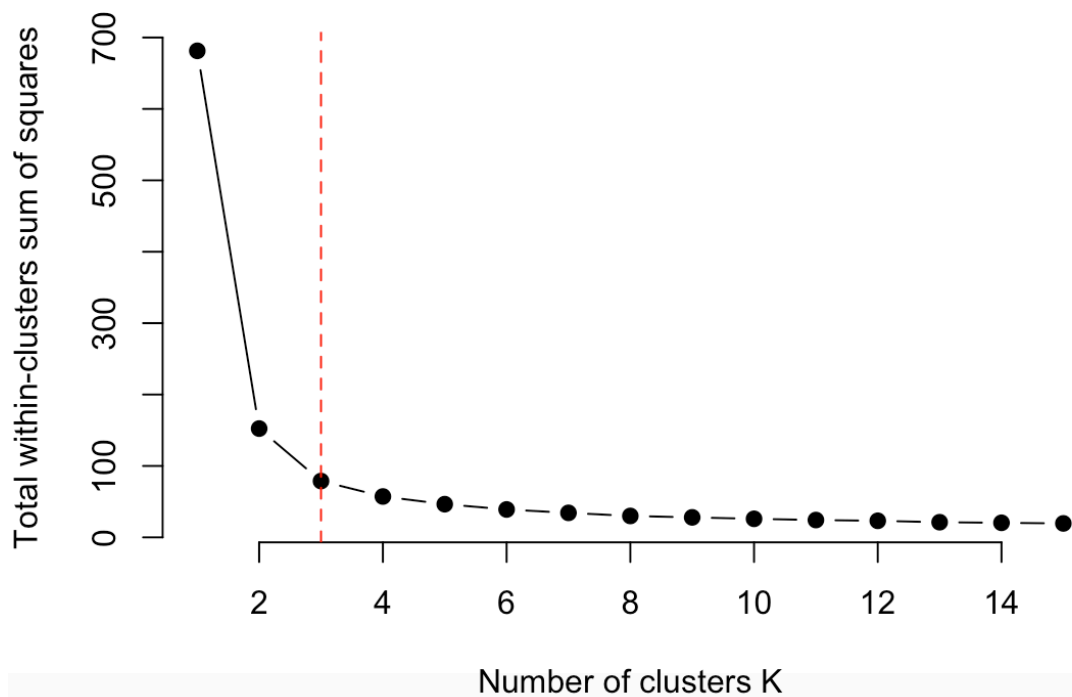
**PCAtemp**

# K-means

```r
k.max <- 15
wss <- sapply(1:k.max, function(k){kmeans(iris2, k, nstart = 10)$tot.withinss})
plot(1:k.max,wss,type = "b",pch = 19, frame = FALSE,
     xlab = "Number of clusters K",ylab = "Total within-clusters sum of squares")
abline(v = 3,lty =2, col = "red")

fviz_nbclust(iris2, kmeans, method = c("silhouette"))

iris.kmeans <- kmeans(iris2,3)
autoplot(iris.kmeans, data = iris2, frame = TRUE, frame.type = "norm")

dis <- dist(iris2)^2
sil <- silhouette(iris.kmeans$cluster, dis)
plot(sil)

plot(iris$Species,iris.kmeans$cluster)
```
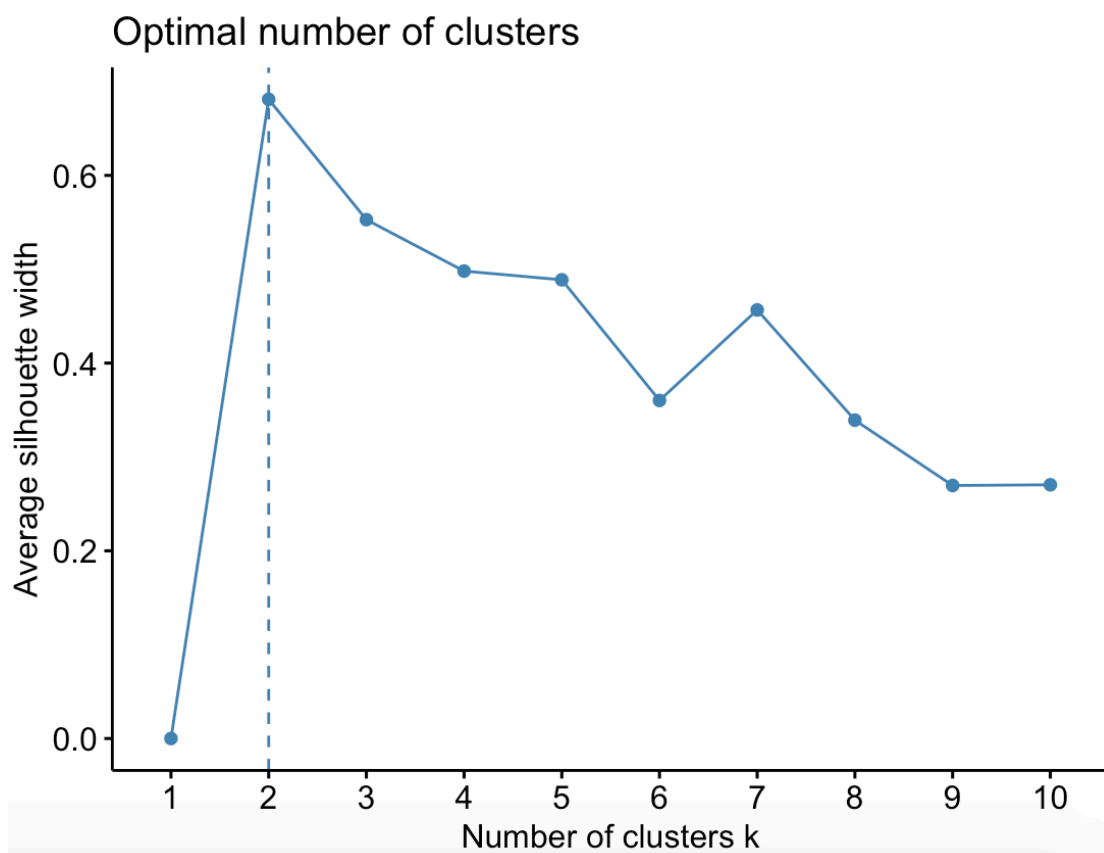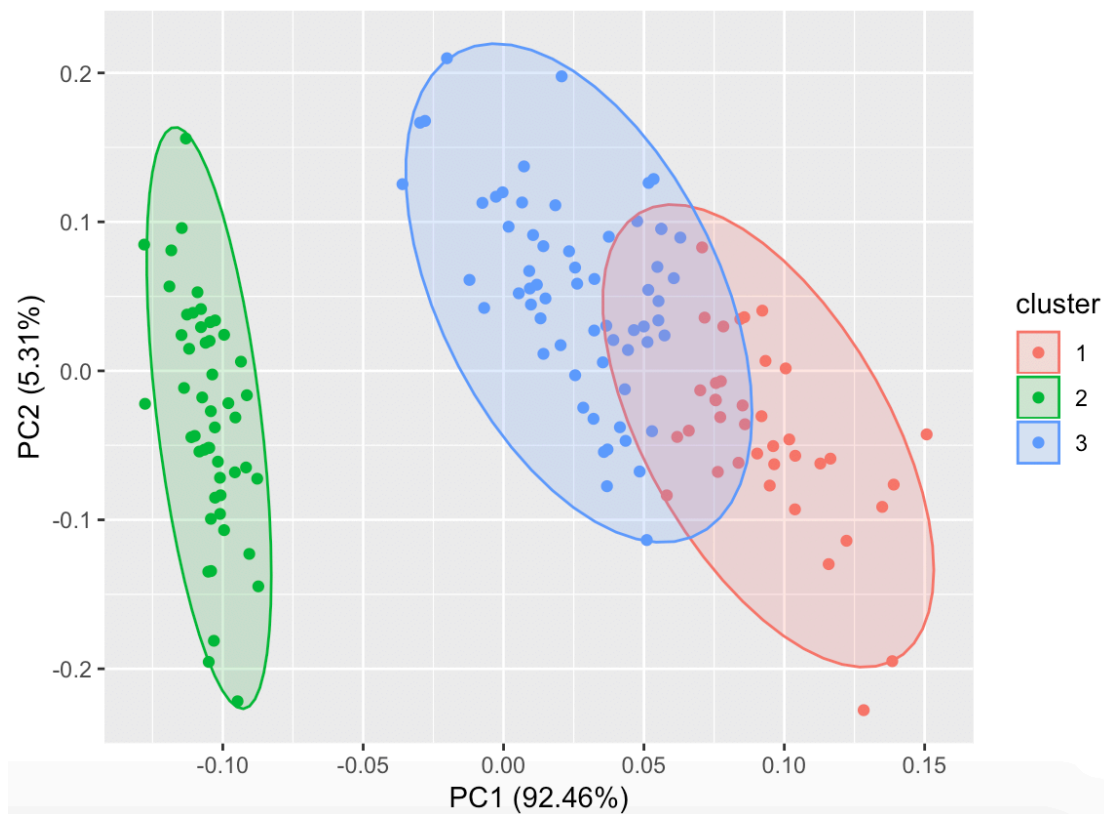
Optimal number of clusters

As shown in the figure, k-means has a good discrimination degree for cluster 1 and a poor discrimination degree for cluster 2 and 3.



**Silhouette plot of (x = iris.kmeans$cluster, dist = dis)**
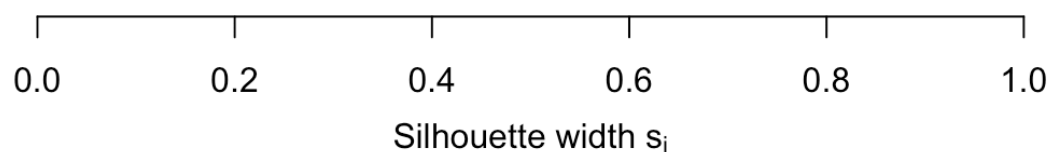
n = 150

3 clusters $C_j$

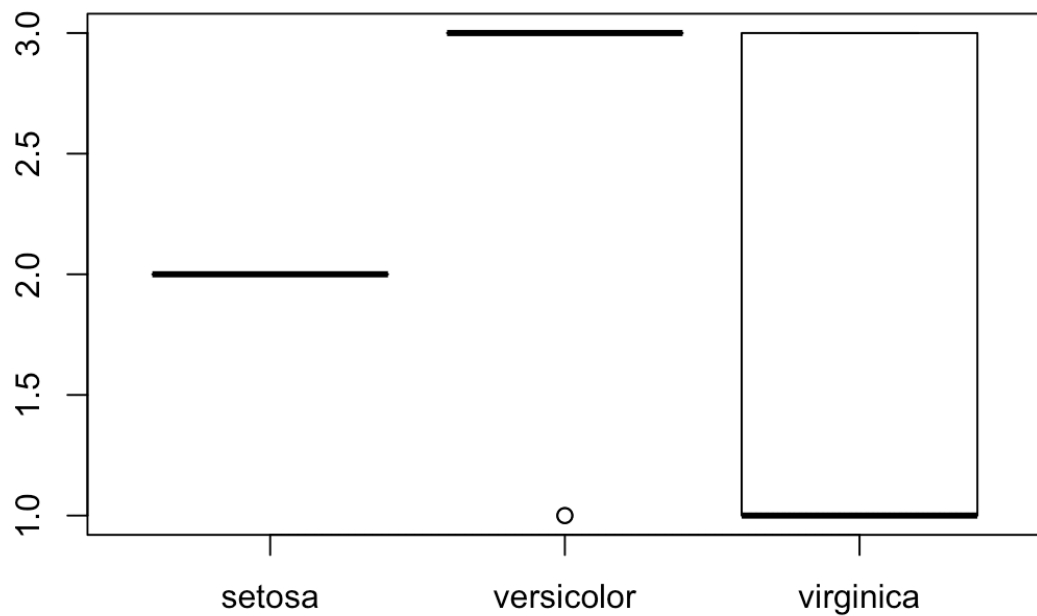$j : n_j$ | $ave_{i \in C_j}$ $s_i$

1 :  38 | 0.66

2 :  50 | 0.95

3 :  62 | 0.61

Silhouette width $s_i$

Average silhouette width :  0.74

```
plot(iris$Species,iris.kmeans$cluster)
```