

PCA etc.

Wolfgang Huber

CSAMA 2015

```

require(ggbiplot)
data(wine)
wine[1:3,1:7]

##   Alcohol MalicAcid  Ash AlcAsh Mg Phenols Flav
## 1    14.2     1.71 2.43 15.6 127   2.80 3.06
## 2    13.2     1.78 2.14 11.2 100   2.65 2.76
## 3    13.2     2.36 2.67 18.6 101   2.80 3.24

```

```
heatmap(1-cor(wine))
```

```
wine.pca <- prcomp(wine, scale. = TRUE)
```

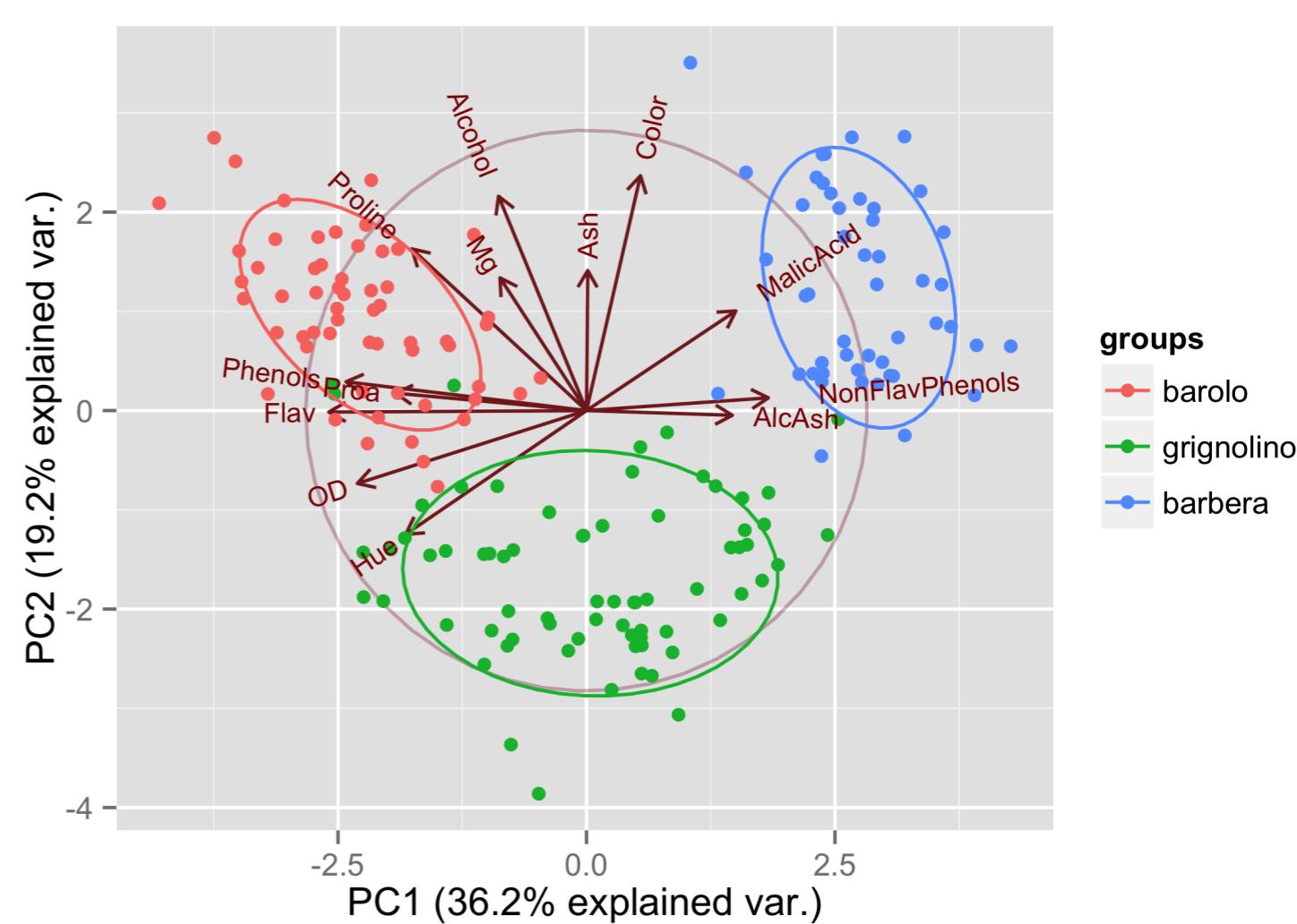
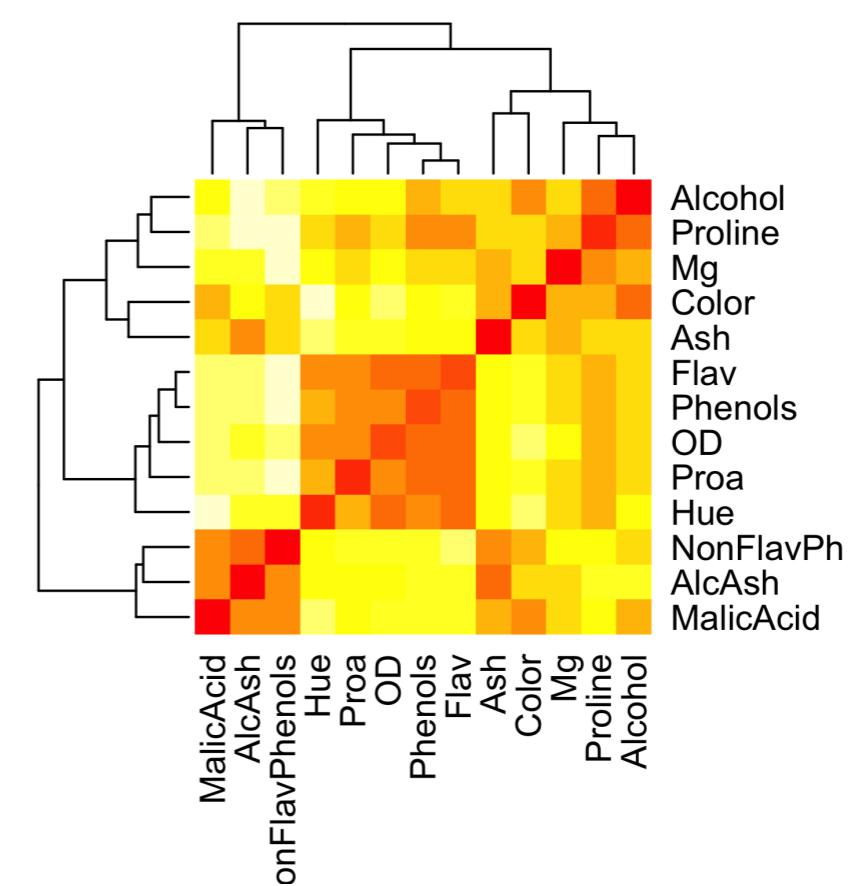
```
table(wine.class)
```

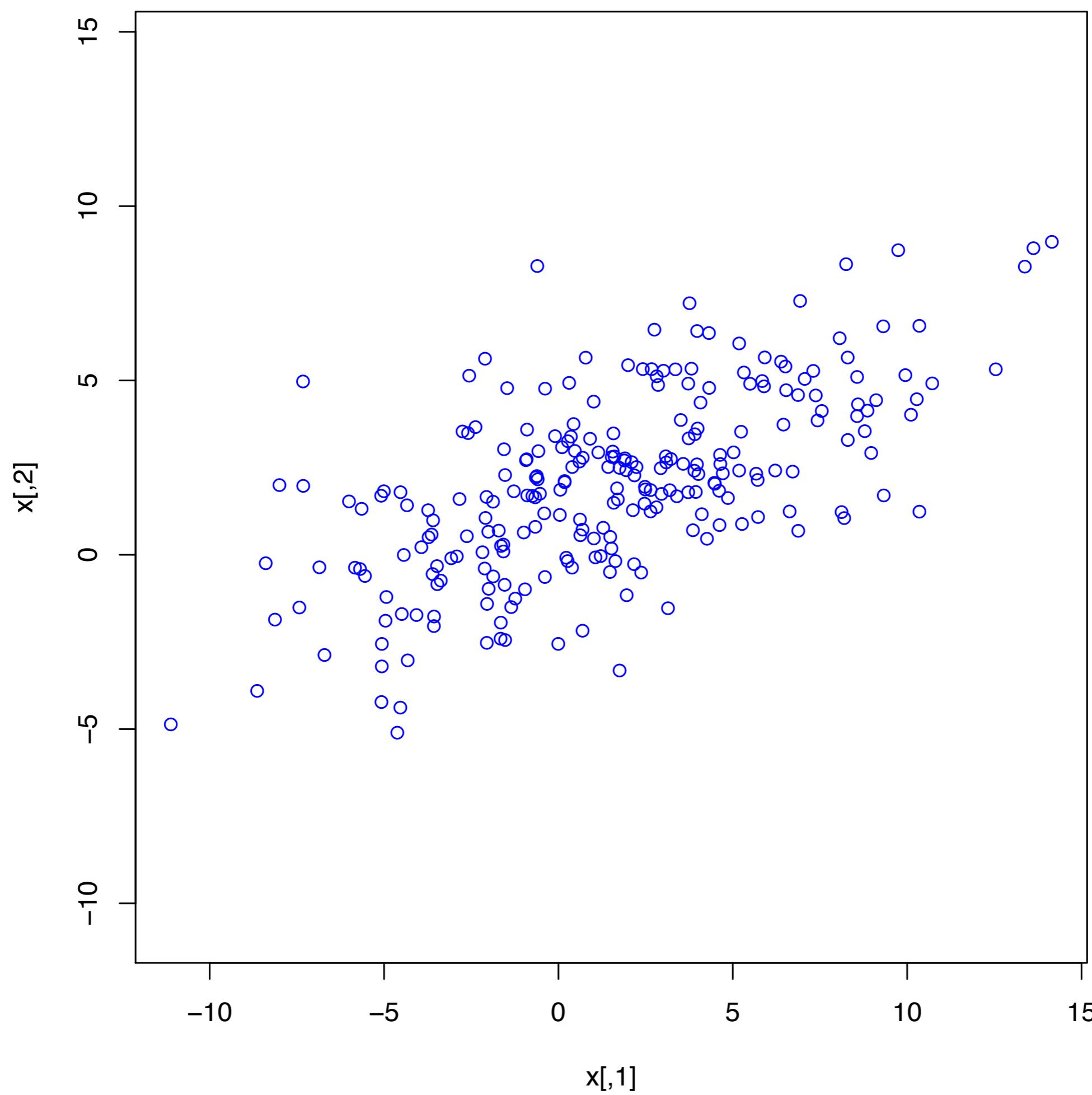
```
## wine.class
##   barolo grignolino  barbera
##      59        71       48
```

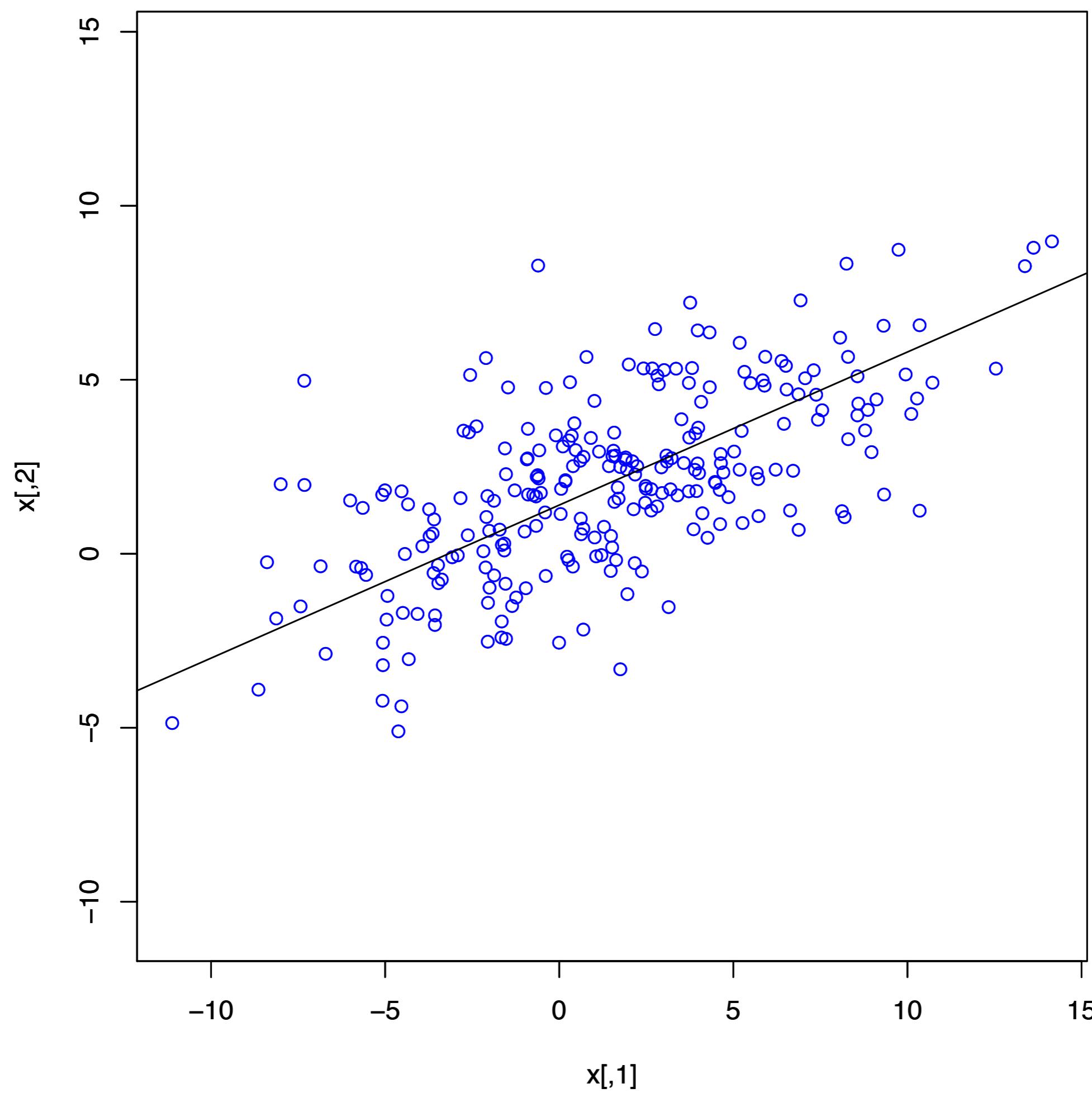
```
pp=ggbiplot(wine.pca, obs.scale = 1, var.scale = 1, scale=1,
```

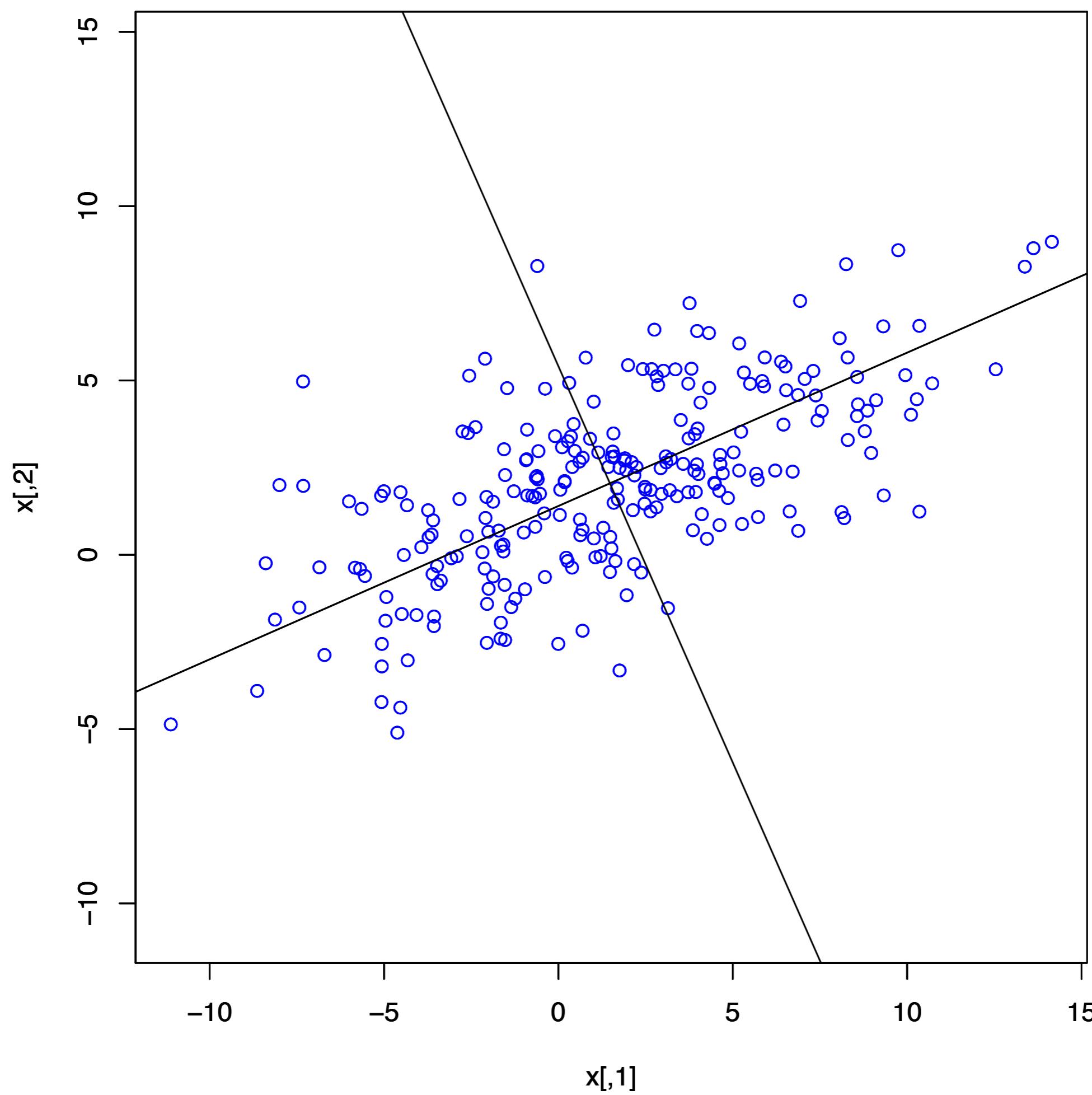
```
  groups = wine.class, ellipse = TRUE, circle = TRUE)
```

```
pp
```









Principal Component Analysis

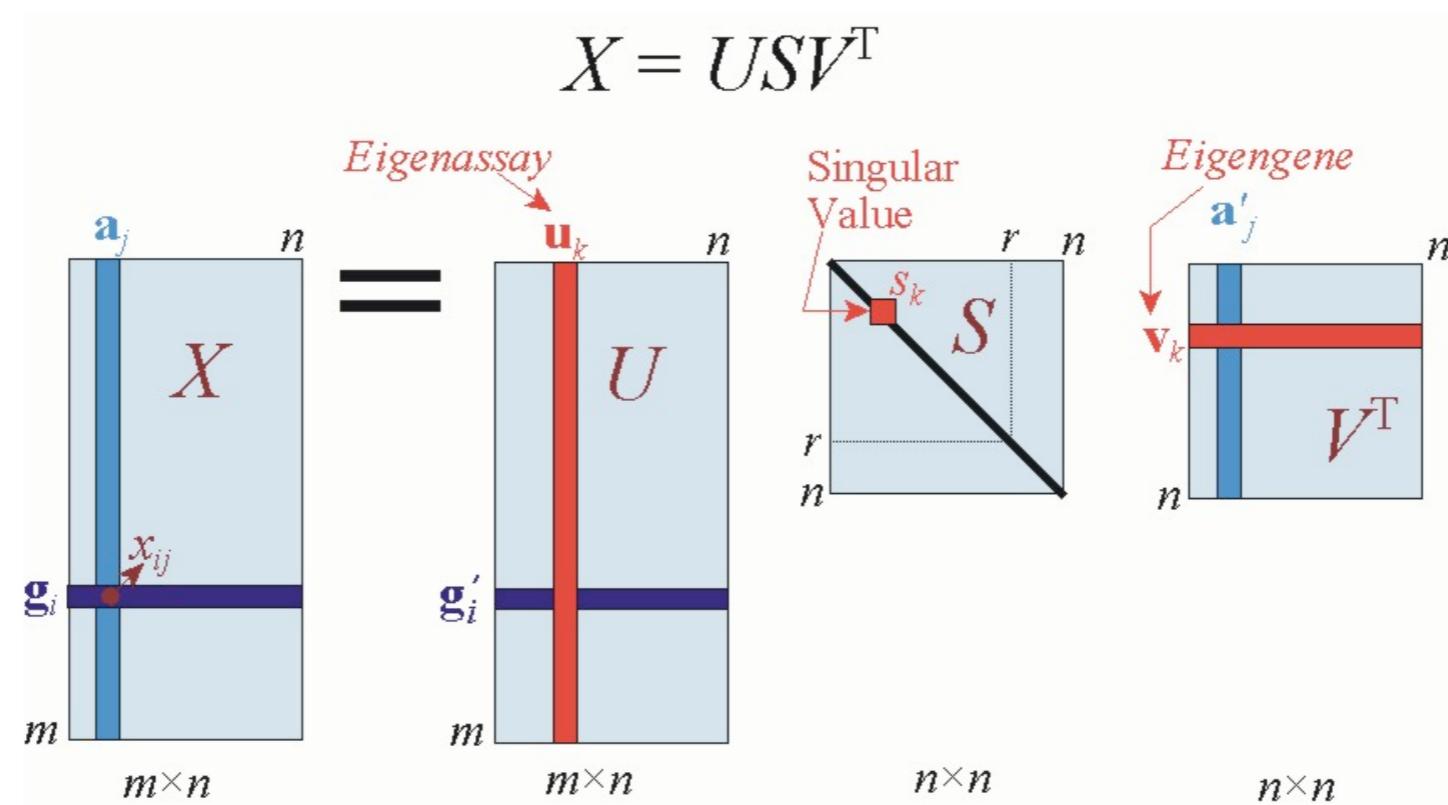
Data points $x_i \in \mathbb{R}^n$

Linear projection $P: \mathbb{R}^n \rightarrow \mathbb{R}^k$ ($n \leq k$)

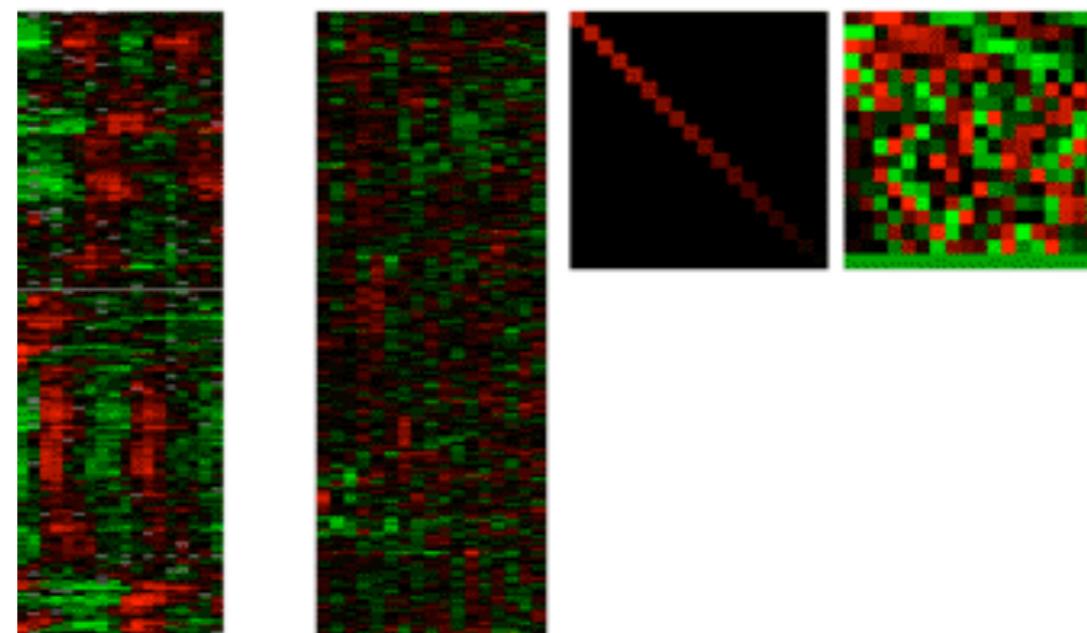
$$\sum_i (x_i - P(x_i))^2 \rightarrow \min$$

$$\text{Cov } P(x_i))^2 \rightarrow \max$$

How is the Principal Component Analysis computed?



$$A = U \cdot W \cdot V^T$$



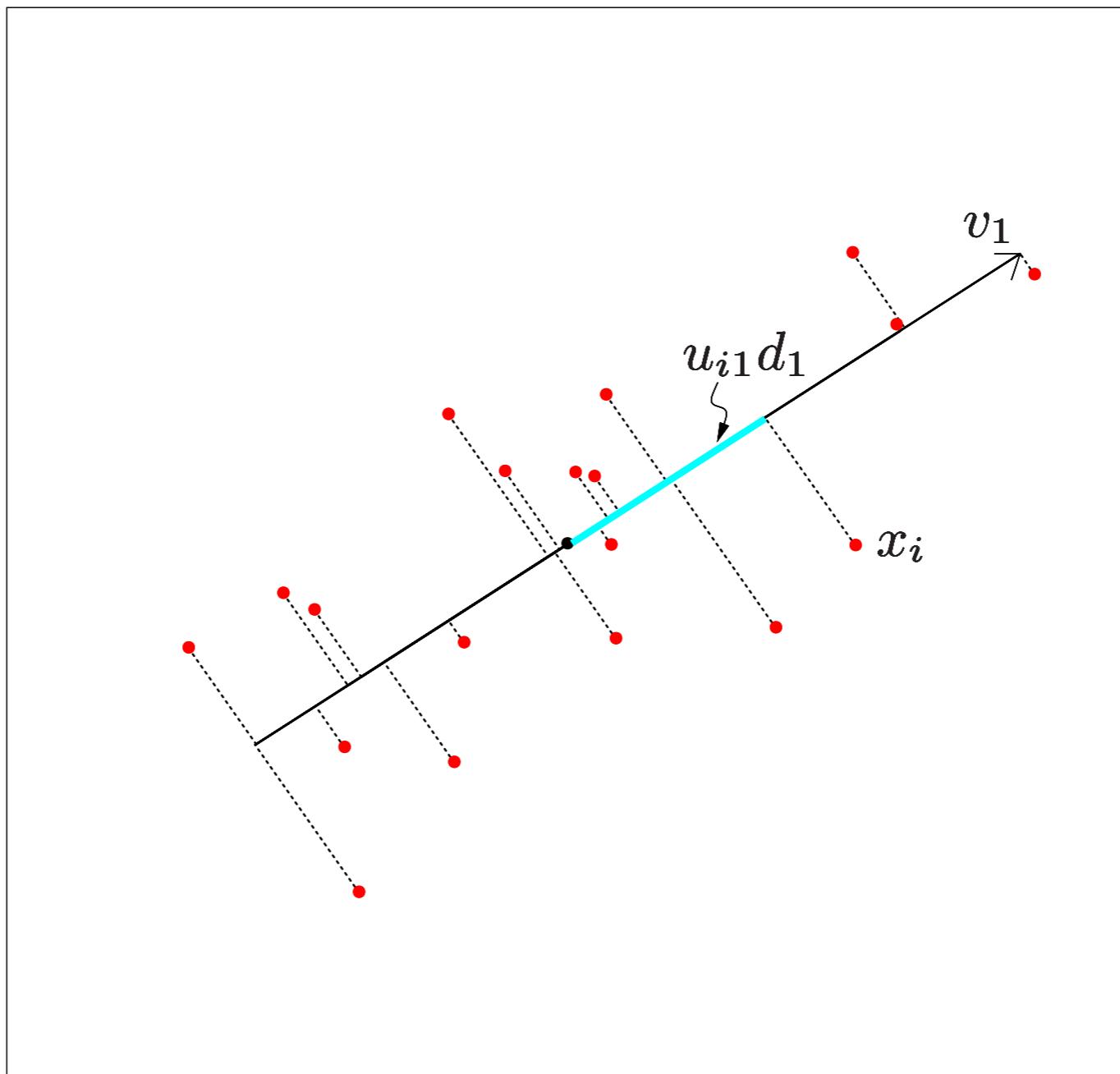


FIGURE 14.20. The first linear principal component of a set of data. The line minimizes the total squared distance from each point to its orthogonal projection onto the line.

Hastie, Tibshirani, Friedman

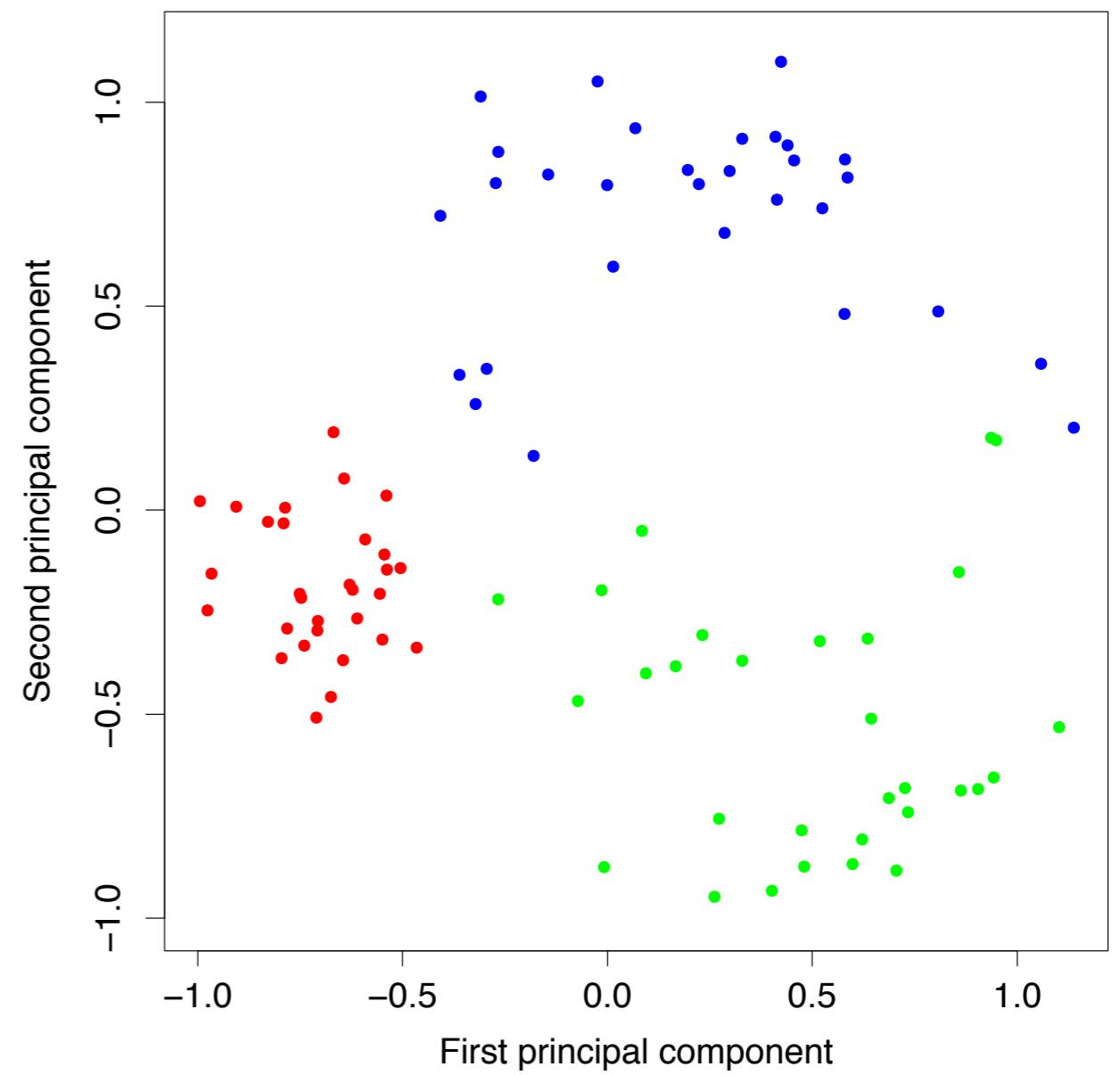
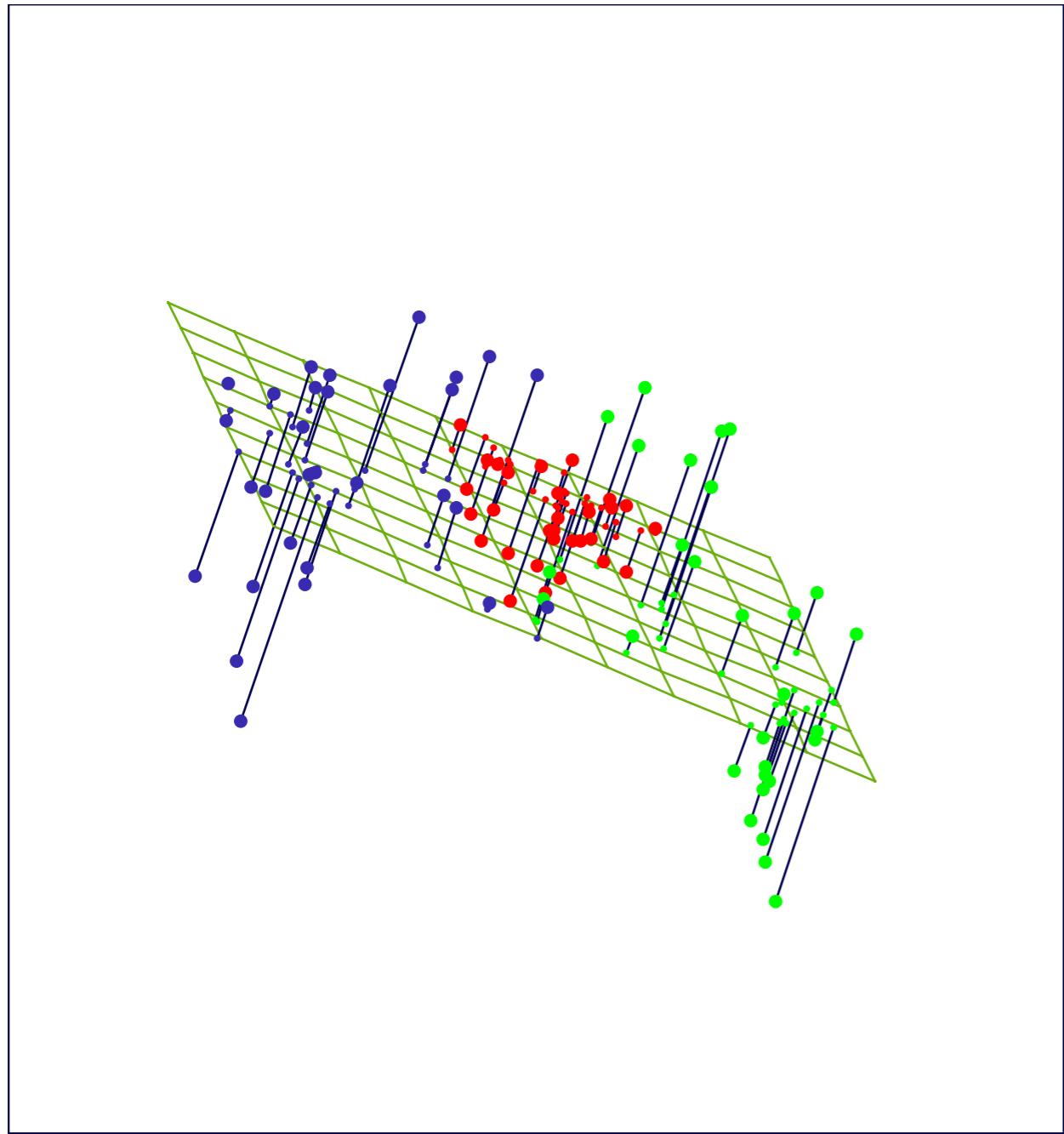


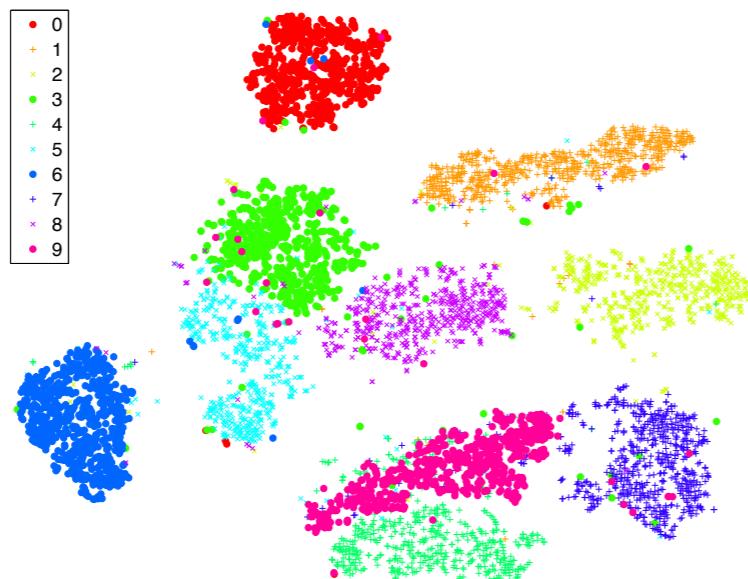
FIGURE 14.21. The best rank-two linear approximation to the half-sphere data. The right panel shows the projected points with coordinates given by $\mathbf{U}_2\mathbf{D}_2$, the first two principal components of the data.

Advanced topics

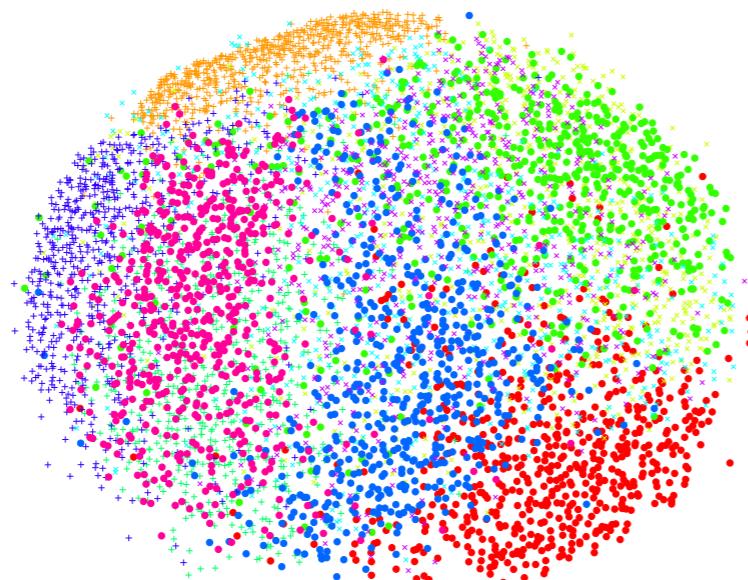
Multi-dimensional scaling: focus on (near) distances

Locally linear embedding (e.g. Donoho & Grimes, PNAS 2003)

t-SNE (van der Maaten & Hinton, JMLR 2008)



(a) Visualization by t-SNE.

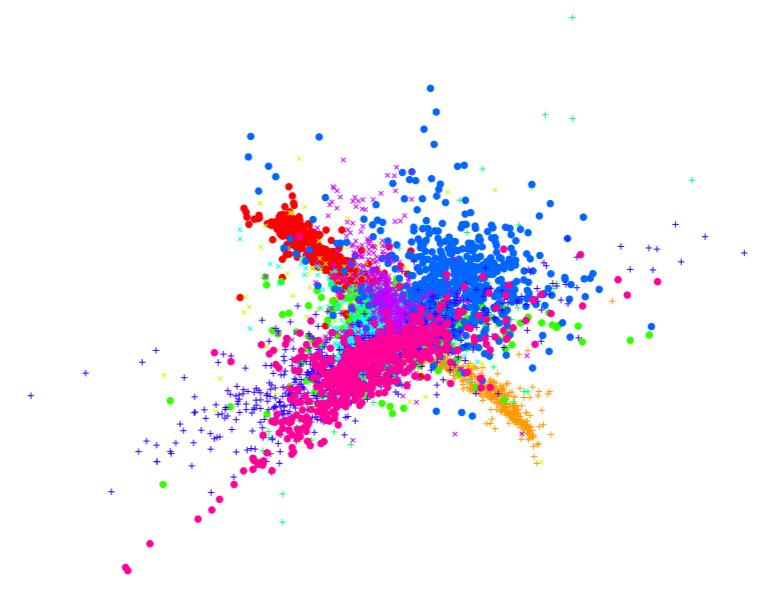


(b) Visualization by Sammon mapping.

Figure 2: Visualizations of 6,000 handwritten digits from the MNIST data set.



(a) Visualization by Isomap.



(b) Visualization by LLE.

Figure 3: Visualizations of 6,000 handwritten digits from the MNIST data set.

Springer Series in Statistics

Trevor Hastie
Robert Tibshirani
Jerome Friedman

The Elements of Statistical Learning

Data Mining, Inference, and Prediction

Second Edition

 Springer

**CRAN Task View: Analysis of Ecological
and Environmental Data:
Ordination**

[http://cran.r-project.org/web/views/
Environmetrics.html](http://cran.r-project.org/web/views/Environmetrics.html)