

Multivariate Statistical Analysis

- 1. Aspects of Multivariate Analysis
- 2. Principal Components
- 3. Factor Analysis
- 4. Discrimination and Classification
- 5. Clustering

Johnson, R.A., Wichern, D.W. (1982): Applied Multivariate Statistical Analysis, Prentice Hall.

1. Aspects of Multivariate Analysis

Multivariate data arise whenever $p \geq 1$ variables are recorded. Values of these variables are observed for n distinct item, individuals, or experimental trials.

We use the notation x_{ij} to indicate the particular value of the i th variable that is observed on the j th item, or trial.

Thus, n measurements on p variables are displayed as $p \times n$ **random** matrix \mathbf{X} :

	Item 1	Item 2	...	Item j	...	Item n
Variable 1:	x_{11}	x_{12}	...	x_{1j}	...	x_{1n}
Variable 2:	x_{21}	x_{22}	...	x_{2j}	...	x_{2n}
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Variable i :	x_{i1}	x_{i2}	...	x_{ij}	...	x_{in}
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Variable p :	x_{p1}	x_{p2}	...	x_{pj}	...	x_{pn}

Estimating Moments:

Suppose, $E(\mathbf{X}) = \boldsymbol{\mu}$ and $\text{cov}(\mathbf{X}) = \boldsymbol{\Sigma}$ are the population moments. Based on a sample of size n , these quantities can be estimated by their empirical versions:

Sample Mean:

$$\bar{x}_i = \frac{1}{n} \sum_{j=1}^n x_{ij}, \quad i = 1, \dots, p$$

Sample Variance:

$$s_i^2 = s_{ii} = \frac{1}{n-1} \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2, \quad i = 1, \dots, p$$

Sample Covariance:

$$s_{ik} = \frac{1}{n-1} \sum_{j=1}^n (x_{ij} - \bar{x}_i)(x_{kj} - \bar{x}_k), \quad i = 1, \dots, p, \quad k = 1, \dots, p.$$

Summarize all elements s_{ik} into the $p \times p$ sample variance-covariance matrix $\mathbf{S} = (s_{ik})_{i,k}$.

Assume further, that the $p \times p$ population correlation matrix $\boldsymbol{\rho}$ is estimated by the sample correlation matrix \mathbf{R} with entries

$$r_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}s_{kk}}}, \quad i = 1, \dots, p, \quad k = 1, \dots, p,$$

where $r_{ii} = 1$ for all $i = 1, \dots, p$.

```
> aimu <- read.table("aimu.dat", header=TRUE)
> attach(aimu)
> options(digits=2)

> mean(aimu[,3:8])
  age height weight   fvc  fev1  fevp
  30   177    77   553   460    83
```

```
> cov(aimu[,3:8])
      age height weight  fvc fev1 fevp
age    110  -16.9   16.5 -233 -302 -20.8
height  -17   45.5   34.9  351  275  -1.9
weight   16   34.9  109.6  325  212  -7.6
fvc    -233  351.5  324.7 5817 4192 -86.5
fev1   -302  275.2  212.0 4192 4347 162.5
fevp    -21  -1.9   -7.6  -87  162  41.3
```

```
> cor(aimu[,3:8])
      age height weight  fvc fev1 fevp
age    1.00 -0.239   0.15 -0.29 -0.44 -0.309
height -0.24  1.000   0.49  0.68  0.62 -0.043
weight  0.15  0.494   1.00  0.41  0.31 -0.113
fvc    -0.29  0.683   0.41  1.00  0.83 -0.177
fev1   -0.44  0.619   0.31  0.83  1.00  0.384
fevp   -0.31 -0.043  -0.11 -0.18  0.38  1.000
```

Distances:

Consider the point $P = (x_1, x_2)$ in the plane. The straight line (Euclidian) distance, $d(O, P)$, from P to the origin $O = (0, 0)$ is (Pythagoras)

$$d(O, P) = \sqrt{x_1^2 + x_2^2}.$$

In general, if P has p coordinates so that $P = (x_1, x_2, \dots, x_p)$, the Euclidian distance is

$$d(O, P) = \sqrt{x_1^2 + x_2^2 + \dots + x_p^2}.$$

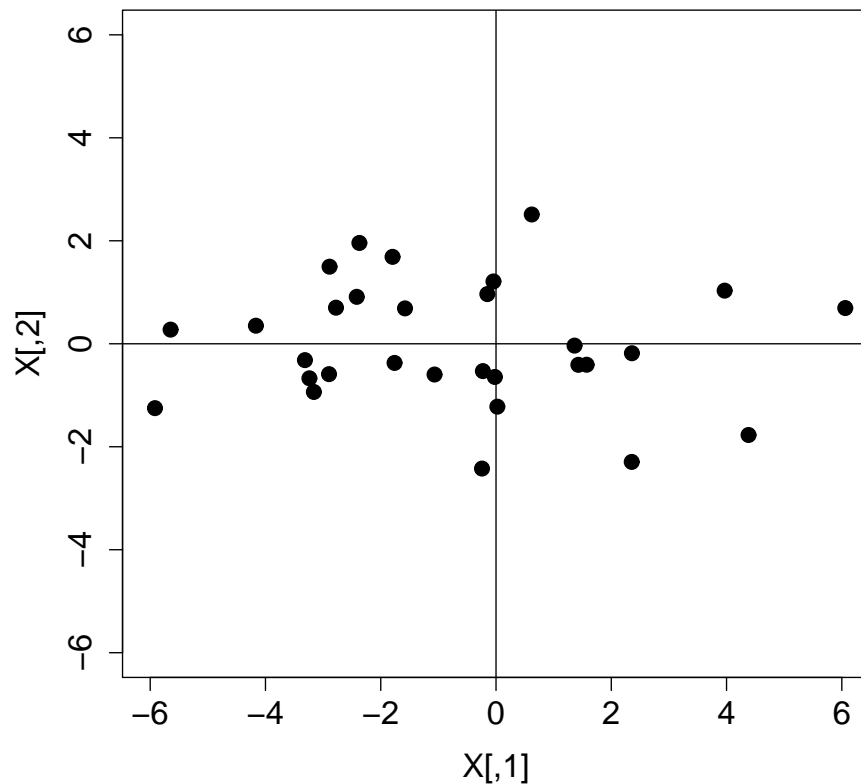
The distance between 2 arbitrary points P and $Q = (y_1, y_2, \dots, y_p)$ is given by

$$d(P, Q) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_p - y_p)^2}.$$

Each coordinate contributes equally to the calculation of the Euclidian distance. It is often desirable to weight the coordinates.

Statistical distance should account for differences in variation and correlation. Suppose we have n pairs of measurements on 2 independent variables x_1 and x_2 :

```
> X <- mvrnorm(30, mu=c(0, 0), Sigma=matrix(c(9,0,0,1), 2, 2)); plot(X)
```



Variability in x_1 direction is much larger than in x_2 direction! Values that are a given deviation from the origin in the x_1 direction are not as *surprising* as are values in x_2 direction.

It seems reasonable to weight an x_2 coordinate more heavily than an x_1 coordinate of the same value when computing the distance to the origin.

Compute the statistical distance from the standardized coordinates

$$x_1^* = \frac{x_1}{\sqrt{s_{11}}} \quad \text{and} \quad x_2^* = \frac{x_2}{\sqrt{s_{22}}}$$

as

$$d(O, P) = \sqrt{(x_1^*)^2 + (x_2^*)^2} = \sqrt{\left(\frac{x_1}{\sqrt{s_{11}}}\right)^2 + \left(\frac{x_2}{\sqrt{s_{22}}}\right)^2} = \sqrt{\frac{x_1^2}{s_{11}} + \frac{x_2^2}{s_{22}}}.$$

This can be generalized to accommodate the calculation of statistical distance from an arbitrary point $P = (x_1, x_2)$ to any *fixed* point $Q = (y_1, y_2)$. If the coordinate variables vary independent of one other, the distance from P to Q is

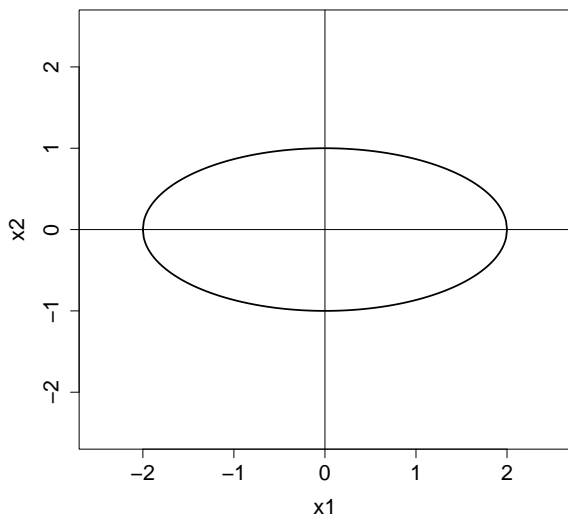
$$d(P, Q) = \sqrt{\frac{(x_1 - y_1)^2}{s_{11}} + \frac{(x_2 - y_2)^2}{s_{22}}}.$$

The extension to more than 2 dimensions is straightforward.

Let $P = (x_1, x_2, \dots, x_p)$ and $Q = (y_1, y_2, \dots, y_p)$. Assume again that Q is fixed. The statistical distance from P to Q is

$$d(P, Q) = \sqrt{\frac{(x_1 - y_1)^2}{s_{11}} + \frac{(x_2 - y_2)^2}{s_{22}} + \dots + \frac{(x_p - y_p)^2}{s_{pp}}}.$$

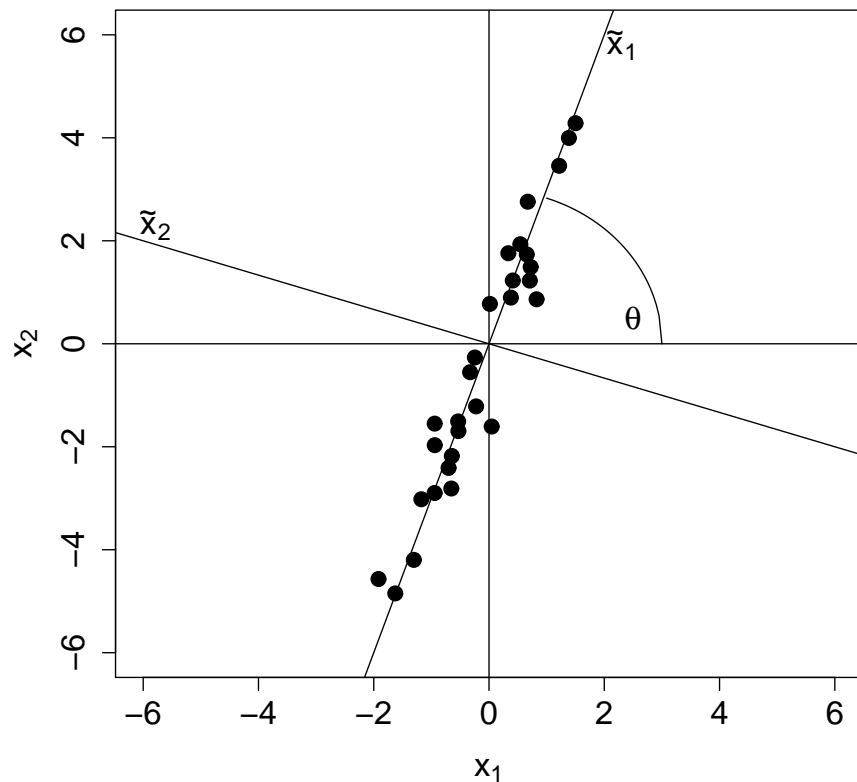
- The distance of P to the origin is obtained by setting $y_1 = y_2 = \dots = y_p = 0$.
- If $s_{11} = s_{22} = \dots = s_{pp}$, the Euclidian distance is appropriate.



Consider a set of paired measurements (x_1, x_2) with $\bar{x}_1 = \bar{x}_2 = 0$, and $s_{11} = 4$, $s_{22} = 1$. Suppose the x_1 measurements are unrelated to the x_2 ones. We measure the squared distance of an arbitrary $P = (x_1, x_2)$ to $(0, 0)$ by $d^2(O, P) = x_1^2/4 + x_2^2/1$. All points with constant distance 1 satisfy: $x_1^2/4 + x_2^2/1 = 1$, an Ellipse centered at $(0, 0)$.

This definition of statistical distance still does not include most of the important cases because of the assumption of independent coordinates.

```
> X <- mvrnorm(30, mu=c(0, 0), Sigma=matrix(c(1,2.9,2.9,9), 2, 2))  
> plot(X); abline(h=0, v=0); abline(0, 3); abline(0, -1/3)
```



Here, the x_1 measurements do not vary independently of x_2 . The coordinates exhibit a tendency to be large or small together. Moreover, the variability in the x_2 directions is larger than in x_1 direction.

What is a meaningful measure of distance? Actually, we can use what we have already introduced!

But before, we only have to rotate the coordinate system through the angle θ and label the rotated axes \tilde{x}_1 and \tilde{x}_2 .

Now, we define the distance of a point $P = (x_1, x_2)$ from the origin $(0, 0)$ as

$$d(O, P) = \sqrt{\frac{\tilde{x}_1^2}{\tilde{s}_{11}} + \frac{\tilde{x}_2^2}{\tilde{s}_{22}}},$$

where \tilde{s}_{ii} denotes the sample variance computed with the (rotated) \tilde{x}_i measurements.

Alternative measures of distance can be useful, provided they satisfy the properties

1. $d(P, Q) = d(Q, P)$,
2. $d(P, Q) > 0$ if $P \neq Q$,
3. $d(P, Q) = 0$ if $P = Q$,
4. $d(P, Q) \leq d(P, R) + d(R, Q)$, R being any other point different to P and Q .

Principle Components (PCA)

Now we try to explain the variance-covariance structure through a few **linear** combinations of the original p variables X_1, X_2, \dots, X_p (data reduction).

Let a random vector $\mathbf{X} = (X_1, X_2, \dots, X_p)^t$ have $p \times p$ population variance-covariance matrix $\text{var}(\mathbf{X}) = \mathbf{\Sigma}$.

Denote the eigenvalues of $\mathbf{\Sigma}$ by $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$.

Consider the arbitrary linear combinations with fixed vectors ℓ_i

$$\begin{aligned} Y_1 &= \ell_1^t \mathbf{X} = \ell_{11}X_1 + \ell_{21}X_2 + \dots + \ell_{p1}X_p \\ Y_2 &= \ell_2^t \mathbf{X} = \ell_{12}X_1 + \ell_{22}X_2 + \dots + \ell_{p2}X_p \\ &\vdots \\ Y_p &= \ell_p^t \mathbf{X} = \ell_{1p}X_1 + \ell_{2p}X_2 + \dots + \ell_{pp}X_p \end{aligned}$$

For these

$$\begin{aligned}\text{var}(Y_i) &= \text{var}(\ell_i^t \mathbf{X}) = \ell_i^t \Sigma \ell_i \\ \text{cov}(Y_i, Y_k) &= \text{cov}(\ell_i^t \mathbf{X}, \ell_k^t \mathbf{X}) = \ell_i^t \Sigma \ell_k\end{aligned}$$

We define as **principal components** those linear combinations Y_1, Y_2, \dots, Y_p , which are **uncorrelated** and whose variances are as **large** as possible.

Since increasing the length of ℓ_i would also increase the variances, we restrict our search onto vectors ℓ_i , which are of unit length, i.e. $\sum_j \ell_{ij}^2 = \ell_i^t \ell_i = 1$.

Procedure:

1. the first principal component is the linear combination $\ell_1^T \mathbf{X}$ that maximizes $\text{var}(\ell_1^t \mathbf{X})$ subject to $\ell_1^t \ell_1 = 1$.
2. the second principal component is the linear combination $\ell_2^T \mathbf{X}$ that maximizes $\text{var}(\ell_2^t \mathbf{X})$ subject to $\ell_2^t \ell_2 = 1$ and with $\text{cov}(\ell_1^t \mathbf{X}, \ell_2^t \mathbf{X}) = 0$ (uncorrelated with the first one).
3. the i th principal component is the linear combination $\ell_i^T \mathbf{X}$ that maximizes $\text{var}(\ell_i^t \mathbf{X})$ subject to $\ell_i^t \ell_i = 1$ and with $\text{cov}(\ell_i^t \mathbf{X}, \ell_k^t \mathbf{X}) = 0$, for $k < i$ (uncorrelated with all the previous ones).

How to find all these vectors ℓ_i ?

We will use well known some results from matrix theory.

Result 1: Let $\text{var}(\mathbf{X}) = \Sigma$ and let Σ have the eigenvalue-eigenvector pairs $(\lambda_1, \mathbf{e}_1), (\lambda_2, \mathbf{e}_2), \dots, (\lambda_p, \mathbf{e}_p)$, where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Then the i th principal component, $i = 1, \dots, p$, is given by

$$Y_i = \mathbf{e}_i^t \mathbf{X} = e_{1i}X_1 + e_{2i}X_2 + \dots + e_{pi}X_p.$$

With this choices

$$\begin{aligned}\text{var}(Y_i) &= \mathbf{e}_i^t \Sigma \mathbf{e}_i = \lambda_i, \\ \text{cov}(Y_i, Y_k) &= \mathbf{e}_i^t \Sigma \mathbf{e}_k = 0.\end{aligned}$$

Thus, the principal components are uncorrelated and have variances equal to the eigenvalues of Σ .

If some λ_i are equal, the choice of the corresponding coefficient vectors \mathbf{e}_i , and hence Y_i , are not unique.

Result 2: Let $Y_1 = \mathbf{e}_1^t \mathbf{X}$, $Y_2 = \mathbf{e}_2^t \mathbf{X}$, ..., $Y_p = \mathbf{e}_p^t \mathbf{X}$ be the principal components. Then

$$\sigma_{11} + \sigma_{22} + \cdots + \sigma_{pp} = \sum_{i=1}^p \text{var}(X_i) = \lambda_1 + \lambda_2 + \cdots + \lambda_p = \sum_{i=1}^p \text{var}(Y_i).$$

Thus, the **total population variance** equals the sum of the eigenvalues. Consequently, the proportion of total variance due to (explained by) the k th principal component is

$$0 < \frac{\lambda_k}{\lambda_1 + \lambda_2 + \cdots + \lambda_p} < 1$$

If most (e.g. 80 to 90%) of the total population variance (for large p) can be attributed to the first one, two, or three principal components, then these components can *replace* the original p variables without much loss of information.

The magnitude of e_{ik} measures the importance of the k th variable to the i th principal component. In particular, e_{ik} is proportional to the correlation coefficient between Y_i and X_k .

Result 3: If $Y_1 = e_1^t \mathbf{X}$, $Y_2 = e_2^t \mathbf{X}$, ..., $Y_p = e_p^t \mathbf{X}$ are the principal components from the variance-covariance matrix Σ , then

$$\rho_{Y_i, X_k} = \frac{e_{ki} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}$$

are the correlation coefficients between the components Y_i and the variables X_k .

It is informative to consider principal components derived from multivariate normal random variables. Suppose $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ having density function

$$f(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-p/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right).$$

Then the $\boldsymbol{\mu}$ centered ellipsoids of constant density are

$$(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = c^2.$$

In the two-dimensional case $\mathbf{x} = (x_1, x_2)^t$ this equals

$$\frac{1}{1 - \rho_{12}^2} \left[\left(\frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}} \right)^2 + \left(\frac{x_2 - \mu_2}{\sqrt{\sigma_{22}}} \right)^2 - 2\rho_{12} \left(\frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}} \right) \left(\frac{x_2 - \mu_2}{\sqrt{\sigma_{22}}} \right) \right] = c^2.$$

These ellipsoids have axes $\pm c\sqrt{\lambda_i} \mathbf{e}_i$, $i = 1, \dots, p$.

Example: Suppose $\mathbf{x} = (x_1, x_2)^t \sim N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with $\boldsymbol{\mu} = (0, 0)^t$ and

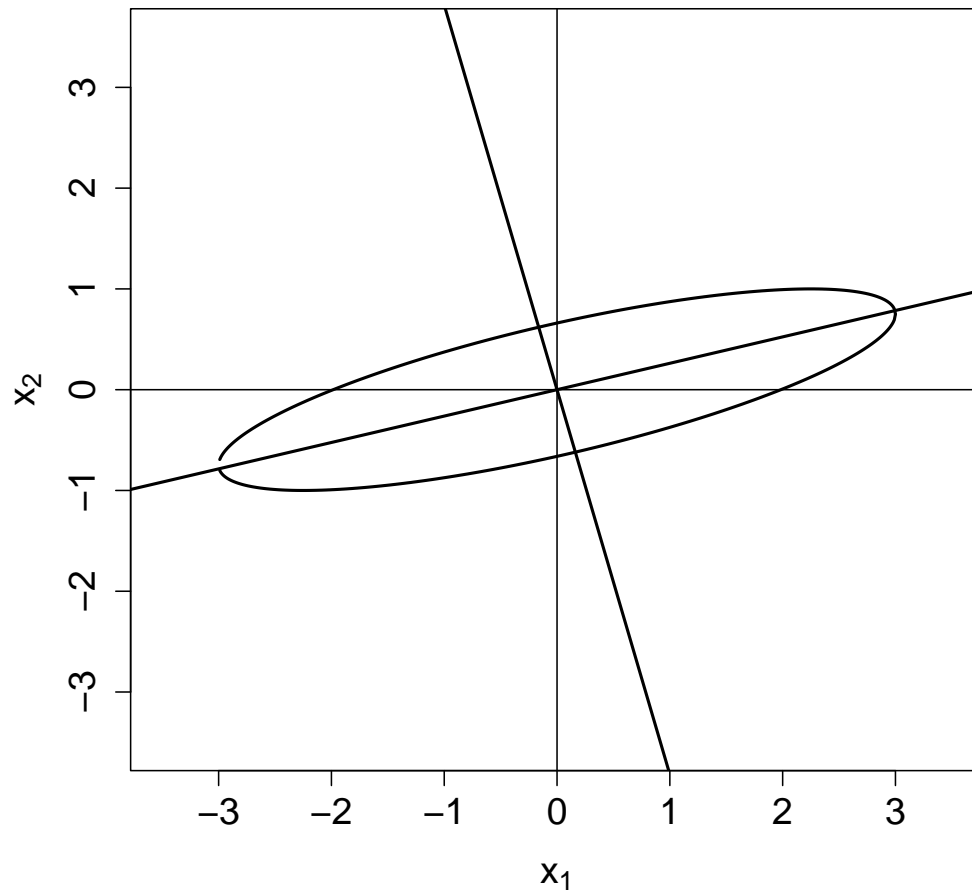
$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{11} = 9 & \sigma_{12} = 9/4 \\ \sigma_{21} = 9/4 & \sigma_{22} = 1 \end{pmatrix}$$

giving $\rho_{12} = (9/4)/\sqrt{9 \cdot 1} = 3/4$.

The eigen-analysis of $\boldsymbol{\Sigma}$ results in

```
> sigma <- matrix(c(9, 9/4, 9/4, 1), 2, 2)
> e <- eigen(sigma, symmetric=TRUE); e
$values
[1] 9.58939 0.41061

$vectors
      [,1]      [,2]
[1,] -0.96736 0.25340
[2,] -0.25340 -0.96736
```



```

# check length of eigenvectors
> e$eigenvectors[2,1]^2+e$eigenvectors[1,1]^2
[1] 1
> e$eigenvectors[2,2]^2+e$eigenvectors[1,2]^2
[1] 1

# slopes of major & minor axes
> e$eigenvectors[2,1]/e$eigenvectors[1,1]
[1] 0.2619511
> e$eigenvectors[2,2]/e$eigenvectors[1,2]
[1] -3.817507

# endpoints of of major&minor axes
> sqrt(e$values[1])*e$eigenvectors[,1]
[1] -2.9956024 -0.7847013
> sqrt(e$values[2])*e$eigenvectors[,2]
[1] 0.1623767 -0.6198741

```

These results also hold for $p \geq 2$. Set $\boldsymbol{\mu} = \mathbf{0}$ in what follows.

$$\begin{aligned}c^2 = \boldsymbol{x}^t \boldsymbol{\Sigma}^{-1} \boldsymbol{x} &= \frac{1}{\lambda_1} (\boldsymbol{e}_1^t \boldsymbol{x})^2 + \frac{1}{\lambda_2} (\boldsymbol{e}_2^t \boldsymbol{x})^2 + \cdots + \frac{1}{\lambda_p} (\boldsymbol{e}_p^t \boldsymbol{x})^2, \\ &= \frac{1}{\lambda_1} y_1^2 + \frac{1}{\lambda_2} y_2^2 + \cdots + \frac{1}{\lambda_p} y_p^2\end{aligned}$$

and this equation defines an ellipsoid (since the λ_i are positive) in a coordinate system with axes y_1, y_2, \dots, y_p lying in the directions of $\boldsymbol{e}_1, \boldsymbol{e}_2, \dots, \boldsymbol{e}_p$. If λ_1 is the largest eigenvalue, then the major axis lies in the direction of \boldsymbol{e}_1 . The remaining minor axes lie in the directions defined by $\boldsymbol{e}_2, \dots, \boldsymbol{e}_p$. Thus the principal components lie in the directions of the axes of the constant density ellipsoid.

Principal Components obtained from Standardized Variables

Instead of using $\mathbf{X} = (X_1, X_2, \dots, X_p)^t$ we now calculate the principal components from $\mathbf{Z} = (Z_1, Z_2, \dots, Z_p)^t$, where

$$Z_i = \frac{X_i - \mu_i}{\sqrt{\sigma_{ii}}}.$$

In matrix notation this equals

$$\mathbf{Z} = \left(\mathbf{V}^{1/2}\right)^{-1} (\mathbf{X} - \boldsymbol{\mu}),$$

where the diagonal standard deviation matrix $\mathbf{V}^{1/2}$ is defined as

$$\mathbf{V}^{1/2} = \begin{pmatrix} \sqrt{\sigma_{11}} & & \\ & \dots & \\ & & \sqrt{\sigma_{pp}} \end{pmatrix}.$$

Clearly $E(\mathbf{Z}) = \mathbf{0}$ and $\text{var}(\mathbf{Z}) = (\mathbf{V}^{1/2})^{-1}\boldsymbol{\Sigma}(\mathbf{V}^{1/2})^{-1} = \boldsymbol{\rho}$.

Principal Components of \mathbf{Z} will be obtained from the eigenvalues λ_i and eigenvectors \mathbf{e}_i of $\boldsymbol{\rho}$ of \mathbf{X} . These are, in general, **not the same** as the ones derived from $\boldsymbol{\Sigma}$.

Result 4: The i th principal component of the standardized variables \mathbf{Z} with $\text{var}(\mathbf{Z}) = \boldsymbol{\rho}$ is given by

$$Y_i = \mathbf{e}_i^t \mathbf{Z} = \mathbf{e}_i^t (\mathbf{V}^{1/2})^{-1} (\mathbf{X} - \boldsymbol{\mu}), \quad i = 1, \dots, p.$$

Moreover,

$$\sum_{i=1}^p \text{var}(Y_i) = \sum_{i=1}^p \text{var}(Z_i) = p.$$

Thus, the proportion explained by the k th principal component is λ_k/p and

$$\rho_{Y_i, Z_k} = \mathbf{e}_{ki} \sqrt{\lambda_i}.$$

Example cont'ed: Let again $\mathbf{x} = (x_1, x_2)^t \sim N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with $\boldsymbol{\mu} = (0, 0)^t$ and

$$\boldsymbol{\Sigma} = \begin{pmatrix} 9 & 9/4 \\ 9/4 & 1 \end{pmatrix} \implies \boldsymbol{\rho} = \begin{pmatrix} 1 & 3/4 \\ 3/4 & 1 \end{pmatrix}.$$

The eigen-analysis of $\boldsymbol{\rho}$ now results in:

```
> rho <- matrix(c(1, 3/4, 3/4, 1), 2, 2)
> e <- eigen(rho, symmetric=TRUE); e
$values
[1] 1.75 0.25

$vectors
      [,1]      [,2]
[1,] 0.70711 0.70711
[2,] 0.70711 -0.70711
```

The total population variance is $p = 2$, and $1.75/2 = 87.5\%$ of this variance is already explained by the first principal component.

The principal components from ρ are

$$Y_1 = 0.707Z_1 + 0.707Z_2 = 0.707\frac{X_1}{3} + 0.707\frac{X_2}{1} = 0.236X_1 + 0.707X_2$$

$$Y_2 = 0.707Z_1 - 0.707Z_2 = 0.707\frac{X_1}{3} - 0.707\frac{X_2}{1} = 0.236X_1 - 0.707X_2,$$

whereas those from Σ have been

$$Y_1 = -0.967X_1 - 0.253X_2$$

$$Y_2 = +0.253X_1 - 0.967X_2.$$

The important first component has explained $9.589/10 = 95.6\%$ of the total variability and is dominated by X_1 (because of its large variance). When the variables are standardized however, the resulting variables contribute equally to the principal components. Variables should be standardized, if they are measured on very different scales.

Summarizing Sample Variation by Principal Components

So far we have dealt with population means μ and variances Σ . If we analyze a sample then we have to replace Σ and μ by their empirical versions S and \bar{x} . The eigenvalues and eigenvectors are then based on S or R instead of Σ or ρ .

```
> library(mva)
> attach(aimu)
> options(digits=2)
> pca <- princomp(aimu[, 3:8])
> summary(pca)
```

Importance of components:

	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5	Comp.6
Standard deviation	96.3	29.443	10.707	7.9581	4.4149	1.30332
Proportion of Variance	0.9	0.084	0.011	0.0061	0.0019	0.00016
Cumulative Proportion	0.9	0.981	0.992	0.9980	0.9998	1.00000

```
> pca$center # the means that were subtracted
  age height weight   fvc   fev1   fevp
  30   177    77   553   460    83
```

```
> pca$scale # the scalings applied to each variable
```

age	height	weight	fvc	fev1	fevp
1	1	1	1	1	1

```
> pca$loadings # a matrix whose columns contain the eigenvectors
```

Loadings:

	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5	Comp.6
age		-0.109	0.645	0.747	0.110	
height			0.119	-0.246	0.960	
weight			0.745	-0.613	-0.251	
fvc	-0.763	-0.624				0.133
fev1	-0.641	0.741				-0.164
fevp		0.212				0.976

	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5	Comp.6
SS loadings	1.00	1.00	1.00	1.00	1.00	1.00
Proportion Var	0.17	0.17	0.17	0.17	0.17	0.17
Cumulative Var	0.17	0.33	0.50	0.67	0.83	1.00

```
> pca$scores # values of the p principal components for each observation
```

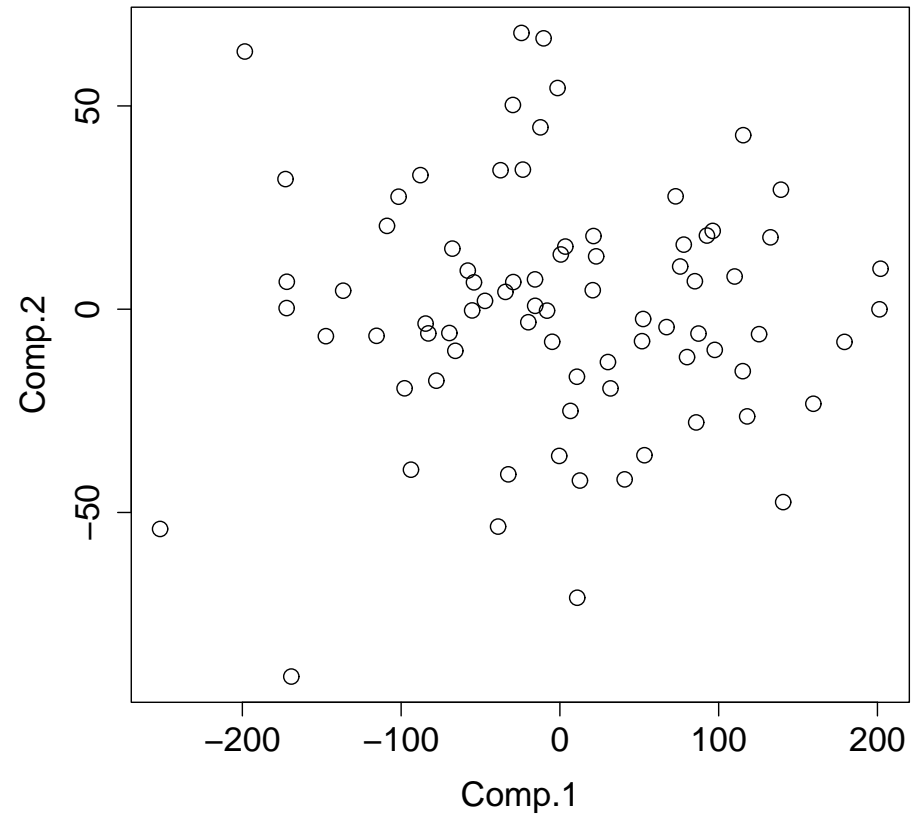
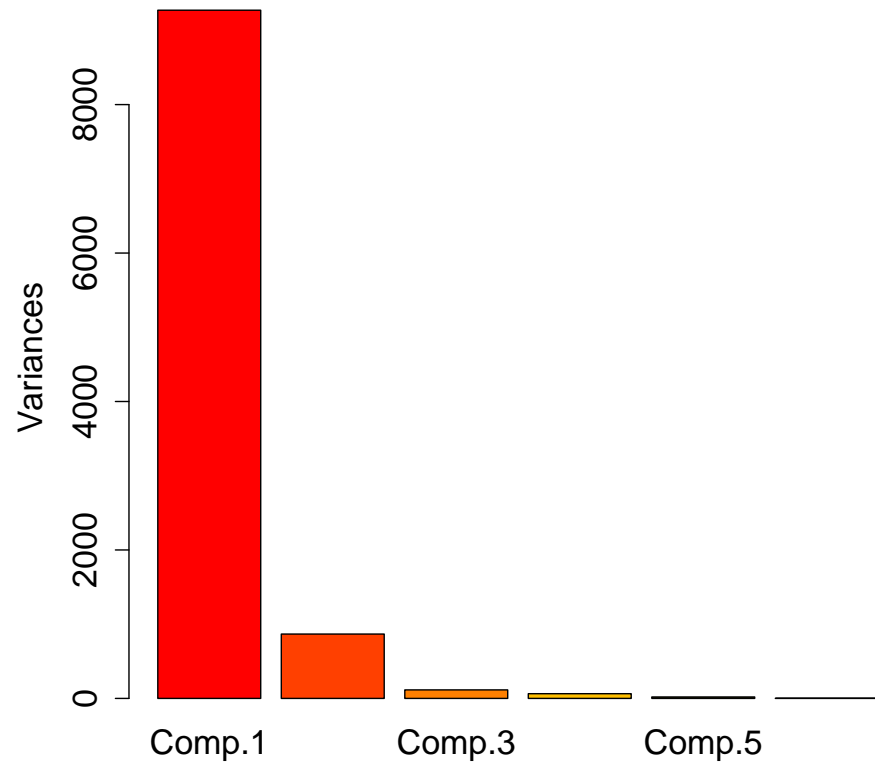
	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5	Comp.6
1	22.84	12.998	4.06	13.131	-1.908	0.0408
2	-147.40	-6.633	-5.14	14.009	-2.130	-0.2862
3	159.64	-23.255	9.60	0.059	5.372	-0.8199
:						
78	52.42	-2.409	1.68	9.169	3.716	0.6386
79	-82.87	-5.951	7.82	11.068	0.834	-0.4171

```
> plot(pca) # or screeplot(pca)
```

```
> plot(pca$scores[, 1:2])
```

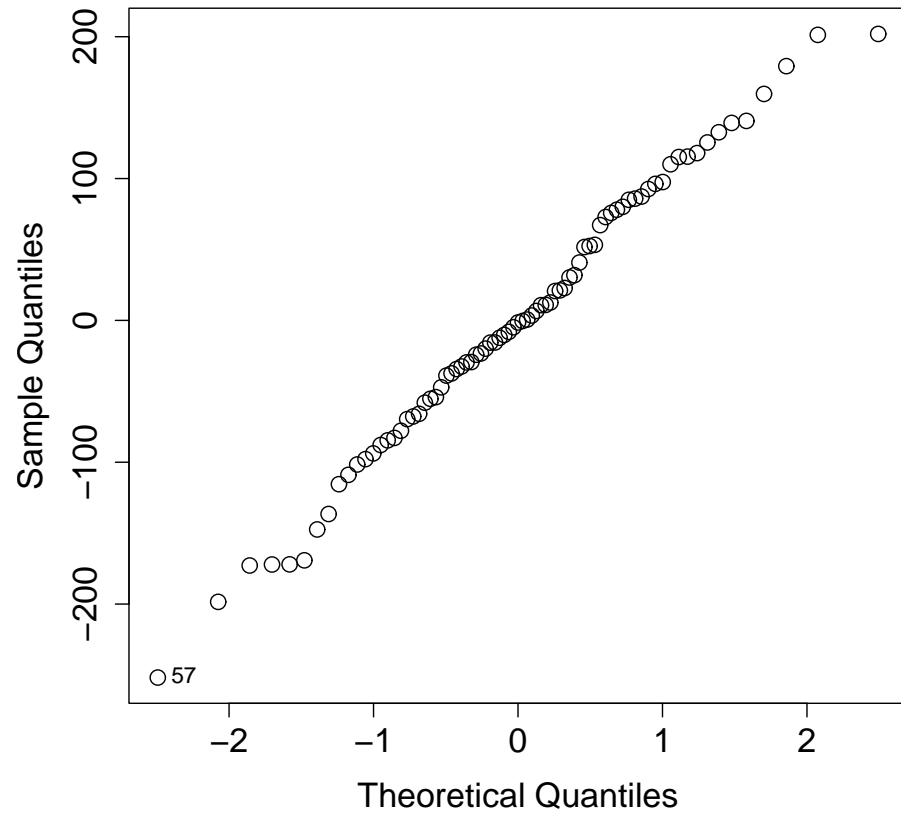
```
> identify(qqnorm(pca$scores[, 1])); identify(qqnorm(pca$scores[, 2]))
```

pca

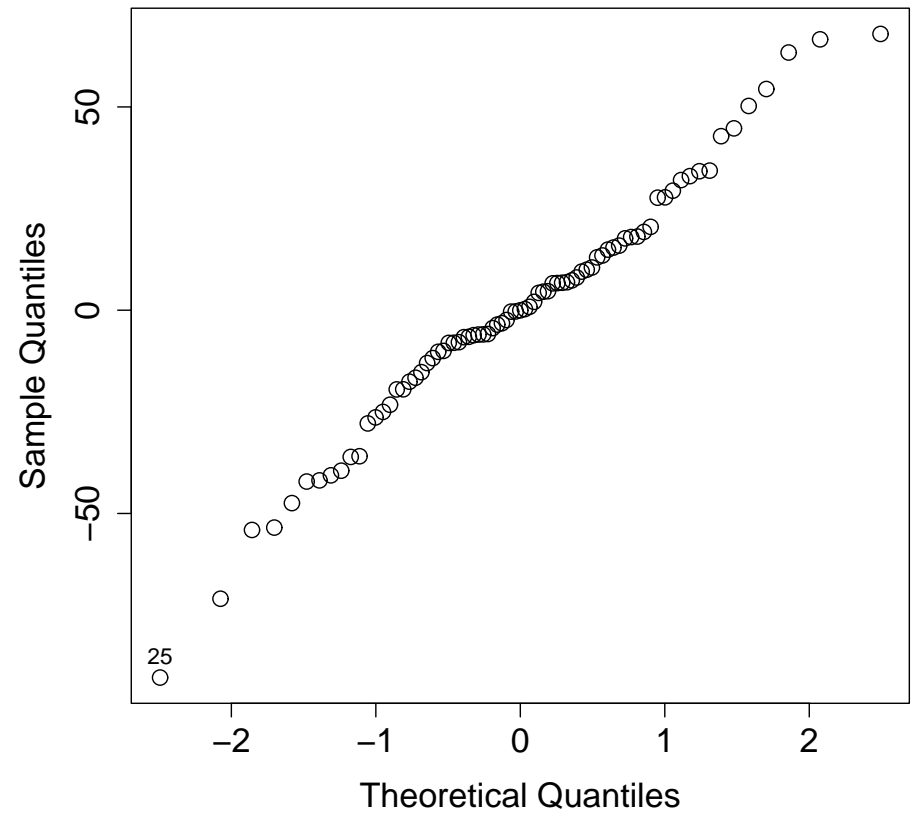


Observations 57 and 25 are a bit outside the ellipsoid.

Normal Q-Q Plot



Normal Q-Q Plot



If we base the analysis on the sample correlation matrix, we get

```
> pca <- princomp(aimu[ , 3:8], cor=TRUE)
```

```
> summary(pca)
```

Importance of components:

	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5	Comp.6
Standard deviation	1.69	1.23	0.91	0.685	0.584	0.0800
Proportion of Variance	0.47	0.25	0.14	0.078	0.057	0.0011
Cumulative Proportion	0.47	0.73	0.86	0.942	0.999	1.0000

```
> pca$center
```

age	height	weight	fv	fev1	fevp
30	177	77	553	460	83

```
> pca$scale
```

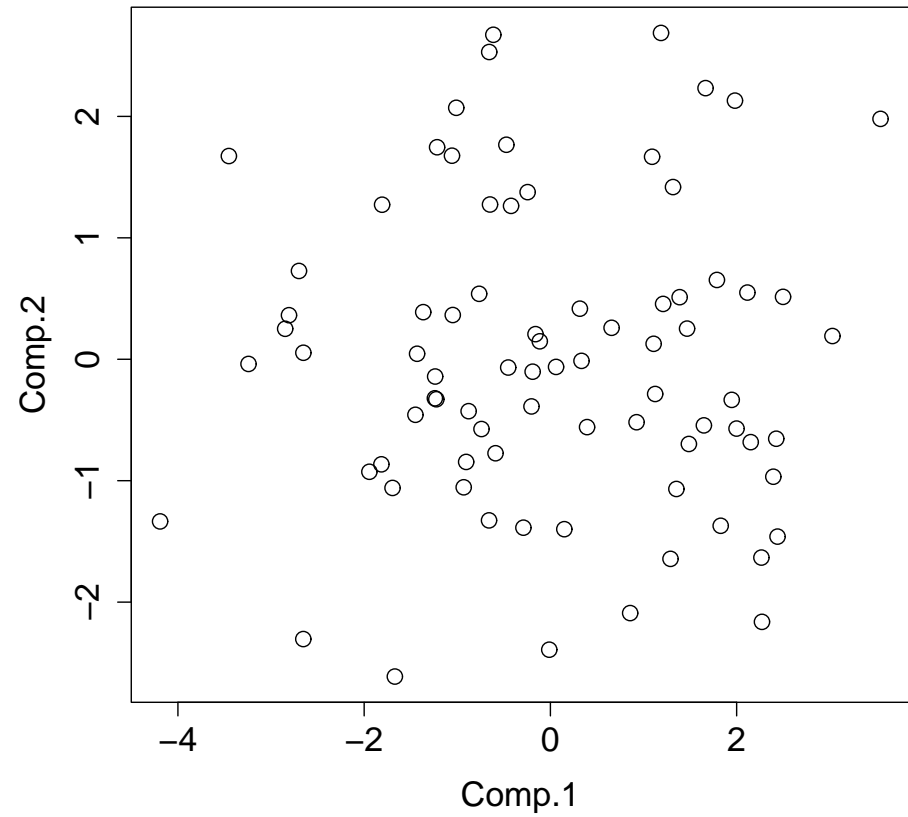
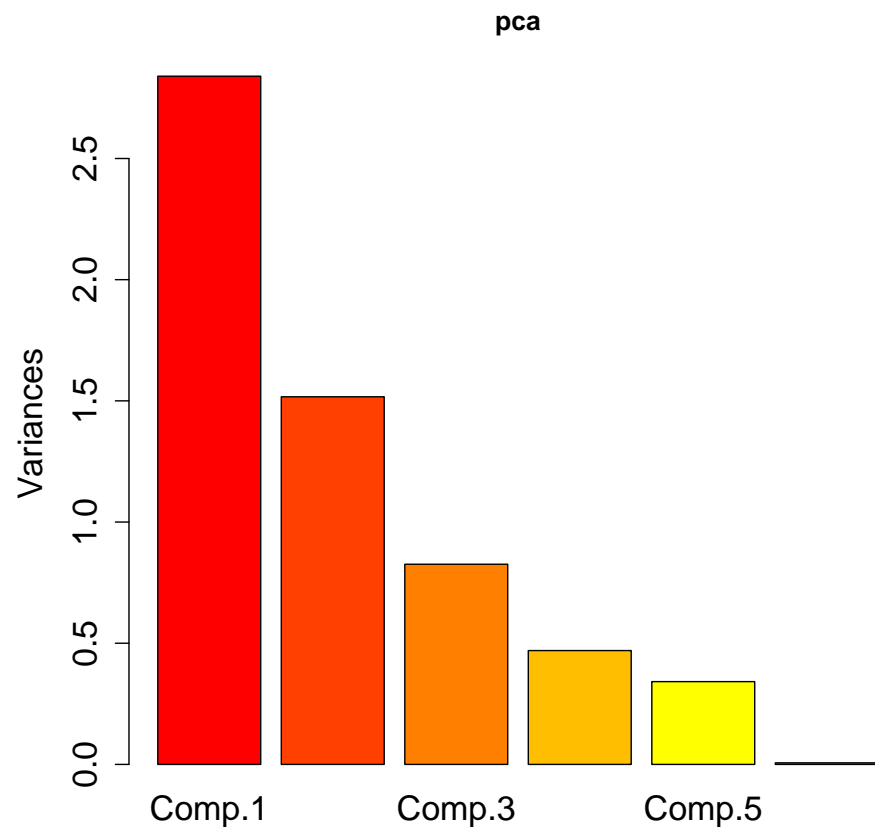
age	height	weight	fv	fev1	fevp
10.4	6.7	10.4	75.8	65.5	6.4

```
> pca$loadings
```

```
Loadings:
```

	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5	Comp.6
age	0.264	-0.535	0.446	0.633	0.211	
height	-0.497	-0.172		-0.207	0.824	
weight	-0.316	-0.449	0.541	-0.494	-0.402	
fvc	-0.534	-0.149	-0.278	0.373	-0.270	0.635
fev1	-0.540	0.217		0.411	-0.168	-0.674
fevp		0.643	0.650		0.110	0.375

	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5	Comp.6
SS loadings	1.00	1.00	1.00	1.00	1.00	1.00
Proportion Var	0.17	0.17	0.17	0.17	0.17	0.17
Cumulative Var	0.17	0.33	0.50	0.67	0.83	1.00



Apart from observations 57 and 25 the plot appears to be reasonable elliptical.

Factor Analysis

Purpose of this (controversial) technique is to describe (if possible) the covariance relationships among many variables in terms of a few underlying but unobservable, random quantities called **factors**.

Suppose variables can be grouped by their correlations. All variables within a group are highly correlated among themselves but have small correlations with variables in a different group. It is conceivable that each such group represents a single underlying construct (factor), that is responsible for the correlations.

E.g., correlations from the group of test scores in French, English, Mathematics suggest an underlying *intelligence factor*. A second group of variables representing *physical fitness scores* might correspond to another factor.

Factor analysis can be considered as an extension of principal component analysis. Both attempt to approximate the covariance matrix Σ .

The Orthogonal Factor Model

The $p \times 1$ random vector \mathbf{X} has mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. The factor model postulates that \mathbf{X} linearly depend on some unobservable random variables F_1, F_2, \dots, F_m , called *common factors* and p additional sources of variation $\epsilon_1, \epsilon_2, \dots, \epsilon_p$, called *errors* or sometimes *specific factors*.

The factor analysis model is

$$\begin{aligned} X_1 - \mu_1 &= \ell_{11}F_1 + \ell_{12}F_2 + \cdots + \ell_{1m}F_m + \epsilon_1 \\ X_2 - \mu_2 &= \ell_{21}F_1 + \ell_{22}F_2 + \cdots + \ell_{2m}F_m + \epsilon_2 \\ &\vdots \\ X_p - \mu_p &= \ell_{p1}F_1 + \ell_{p2}F_2 + \cdots + \ell_{pm}F_m + \epsilon_p \end{aligned}$$

or in matrix notation

$$\underbrace{\mathbf{X} - \boldsymbol{\mu}}_{(p \times 1)} = \underbrace{\mathbf{L}}_{(p \times m)} \underbrace{\mathbf{F}}_{(m \times 1)} + \underbrace{\boldsymbol{\epsilon}}_{(p \times 1)} .$$

The coefficient ℓ_{ij} is called *loading* of the i th variable on the j th factor, so \mathbf{L} is the matrix of factor loadings. Notice, that the p deviations $X_i - \mu_i$ are expressed in terms of $p + m$ random variables F_1, \dots, F_m and $\epsilon_1, \dots, \epsilon_p$, which are all **unobservable**. (This distinguishes the factor model from a regression model, where the explanatory variables F_j can be observed.)

There are too many unobservable quantities in the model. Hence we need further assumptions about \mathbf{F} and ϵ . We assume that

$$\begin{aligned} \mathbf{E}(\mathbf{F}) &= \mathbf{0}, & \text{cov}(\mathbf{F}) &= \mathbf{E}(\mathbf{F}\mathbf{F}^t) = \mathbf{I} \\ \mathbf{E}(\epsilon) &= \mathbf{0}, & \text{cov}(\epsilon) &= \mathbf{E}(\epsilon\epsilon^t) = \boldsymbol{\psi} = \begin{pmatrix} \psi_1 & 0 & \dots & 0 \\ 0 & \psi_2 & \dots & 0 \\ 0 & 0 & \dots & \psi_p \end{pmatrix}. \end{aligned}$$

and \mathbf{F} and ϵ are independent, so

$$\text{cov}(\epsilon, \mathbf{F}) = \mathbf{E}(\epsilon\mathbf{F}^t) = \mathbf{0}.$$

This defines the **orthogonal factor model** and implies a covariance structure for \mathbf{X} . Because of

$$\begin{aligned}
 (\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^t &= (\mathbf{LF} + \boldsymbol{\epsilon})(\mathbf{LF} + \boldsymbol{\epsilon})^t \\
 &= (\mathbf{LF} + \boldsymbol{\epsilon})((\mathbf{LF})^t + \boldsymbol{\epsilon}^t) \\
 &= \mathbf{LF}(\mathbf{LF})^t + \boldsymbol{\epsilon}(\mathbf{LF})^t + (\mathbf{LF})\boldsymbol{\epsilon}^t + \boldsymbol{\epsilon}\boldsymbol{\epsilon}^t
 \end{aligned}$$

we have

$$\begin{aligned}
 \boldsymbol{\Sigma} &= \text{cov}(\mathbf{X}) = \text{E}((\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^t) \\
 &= \mathbf{L} \text{E}(\mathbf{F}\mathbf{F}^t) \mathbf{L}^t + \text{E}(\boldsymbol{\epsilon}\mathbf{F}^t) \mathbf{L}^t + \mathbf{L} \text{E}(\mathbf{F}\boldsymbol{\epsilon}^t) + \text{E}(\boldsymbol{\epsilon}\boldsymbol{\epsilon}^t) \\
 &= \mathbf{L}\mathbf{L}^t + \boldsymbol{\psi}.
 \end{aligned}$$

Since $(\mathbf{X} - \boldsymbol{\mu})\mathbf{F}^t = (\mathbf{LF} + \boldsymbol{\epsilon})\mathbf{F}^t = \mathbf{L}\mathbf{F}\mathbf{F}^t + \boldsymbol{\epsilon}\mathbf{F}^t$ we further get

$$\text{cov}(\mathbf{X}, \mathbf{F}) = \text{E}((\mathbf{X} - \boldsymbol{\mu})\mathbf{F}^t) = \text{E}(\mathbf{L}\mathbf{F}\mathbf{F}^t + \boldsymbol{\epsilon}\mathbf{F}^t) = \mathbf{L} \text{E}(\mathbf{F}\mathbf{F}^t) + \text{E}(\boldsymbol{\epsilon}\mathbf{F}^t) = \mathbf{L}.$$

That proportion of $\text{var}(X_i) = \sigma_{ii}$ contributed by the m common factors is called the i th **communality** h_i^2 . The proportion of $\text{var}(X_i)$ due to the specific factor is called the **uniqueness**, or **specific variance**. I.e.,

$$\text{var}(X_i) = \text{communality} + \text{specific variance}$$

$$\sigma_{ii} = \ell_{i1}^2 + \ell_{i2}^2 + \cdots + \ell_{im}^2 + \psi_i.$$

With $h_i^2 = \ell_{i1}^2 + \ell_{i2}^2 + \cdots + \ell_{im}^2$ we get

$$\sigma_{ii}^2 = h_i^2 + \psi_i.$$

The factor model assumes that the $p(p+1)/2$ variances and covariances of \mathbf{X} can be reproduced by the pm factor loadings ℓ_{ij} and the p specific variances ψ_i . For $p = m$, the matrix Σ can be reproduced exactly as $\mathbf{L}\mathbf{L}^t$, so ψ is the zero matrix. If m is small relative to p , then the factor model provides a simple explanation of Σ with fewer parameters.

Drawbacks:

- Most covariance matrices can not be factored as $\mathbf{L}\mathbf{L}^t + \psi$, where $m \ll p$.
- There is some inherent ambiguity associated with the factor model: let \mathbf{T} be any $m \times m$ orthogonal matrix so that $\mathbf{T}\mathbf{T}^t = \mathbf{T}\mathbf{T} = \mathbf{I}$. then we can rewrite the factor model as

$$\mathbf{X} - \boldsymbol{\mu} = \mathbf{L}\mathbf{F} + \boldsymbol{\epsilon} = \mathbf{L}\mathbf{T}\mathbf{T}^t\mathbf{F} + \boldsymbol{\epsilon} = \mathbf{L}^*\mathbf{F}^* + \boldsymbol{\epsilon}.$$

Since with $\mathbf{L}^* = \mathbf{L}\mathbf{T}$ and $\mathbf{F}^* = \mathbf{T}\mathbf{F}$ we also have

$$\mathbf{E}(\mathbf{F}^*) = \mathbf{T}\mathbf{E}(\mathbf{F}) = \mathbf{0}, \quad \text{and} \quad \text{cov}(\mathbf{F}^*) = \mathbf{T}^t \text{cov}(\mathbf{F})\mathbf{T} = \mathbf{T}^t\mathbf{T} = \mathbf{I},$$

it is impossible to distinguish the loadings in \mathbf{L} from those in \mathbf{L}^* . The factors \mathbf{F} and \mathbf{F}^* have the same statistical properties.

Methods of Estimation

With observations $\boldsymbol{x}_1, \boldsymbol{x}_2, \dots, \boldsymbol{x}_n$ on \boldsymbol{X} , factor analysis seeks to answer the question: Does the factor model with a smaller number of factors adequately represent the data?

The sample covariance matrix \boldsymbol{S} is an estimator of the unknown population covariance matrix $\boldsymbol{\Sigma}$. If the off-diagonal elements of \boldsymbol{S} are small, the variables are not related and a factor analysis model will not prove useful. In these cases, the *specific variances* play the dominant role, whereas the major aim of factor analysis is to determine a few important *common factors*.

If \boldsymbol{S} deviate from a diagonal matrix then the initial problem is to estimate the factor loadings \boldsymbol{L} and specific variances $\boldsymbol{\psi}$. Two methods are very popular: the *principal component method* and the *maximum likelihood method*. Both of these solutions can be *rotated* afterwards in order to simplify the interpretation of the factors.

The Principal Component Approach:

Let Σ have eigenvalue-eigenvector pairs $(\lambda_i, \mathbf{e}_i)$ with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$. Then

$$\Sigma = \lambda_1 \mathbf{e}_1 \mathbf{e}_1^t + \lambda_2 \mathbf{e}_2 \mathbf{e}_2^t + \dots + \lambda_p \mathbf{e}_p \mathbf{e}_p^t.$$

Thus we define

$$\mathbf{L}^t = (\sqrt{\lambda_1} \mathbf{e}_1, \sqrt{\lambda_2} \mathbf{e}_2, \dots, \sqrt{\lambda_p} \mathbf{e}_p)$$

to get a factor analysis model with as many factors as variables ($m = p$) and specific variances $\psi_i = 0$ for all i i.e.

$$\Sigma = \mathbf{L}\mathbf{L}^t + \mathbf{0} = \mathbf{L}\mathbf{L}^t.$$

This is not very useful, however, if the last eigenvalues are relatively small we neglect the contributions of $\lambda_{m+1} \mathbf{e}_{m+1} \mathbf{e}_{m+1}^t + \lambda_{m+2} \mathbf{e}_{m+2} \mathbf{e}_{m+2}^t + \dots + \lambda_p \mathbf{e}_p \mathbf{e}_p^t$ to Σ above.

This gives us the approximation

$$\Sigma \approx \lambda_1 \mathbf{e}_1 \mathbf{e}_1^t + \lambda_2 \mathbf{e}_2 \mathbf{e}_2^t + \cdots + \lambda_m \mathbf{e}_m \mathbf{e}_m^t = \mathbf{L} \mathbf{L}^t,$$

where \mathbf{L} is now a $(m \times p)$ matrix of coefficients as required. This representation assumes that the specific factors ϵ are of minor importance. If specific factors are included in the model, their variances may be taken to be the diagonal elements of $\Sigma - \mathbf{L} \mathbf{L}^t$ and the approximation becomes

$$\Sigma \approx \mathbf{L} \mathbf{L}^t + \psi,$$

where $\psi_i = \sigma_{ii}^2 - \sum_j \ell_{ij}^2$.

To apply this approach to data, it is customary first to center the observations (this does not change the sample covariance structure) and to consider $\mathbf{x}_j - \bar{\mathbf{x}}$.

If the units of the variables are not of the same size then it is desirable to work with the standardized variables $z_{ij} = (x_{ij} - \bar{x}_i) / \sqrt{s_{ii}}$ having sample variance \mathbf{R} .

Applying the above technique onto \mathbf{S} or \mathbf{R} is known as the *principal component solution*.

By the definition of $\hat{\psi}_i = s_{ii} - \sum_j \hat{\ell}_{ij}^2$, where $\hat{\ell}_i$ are the eigenvectors of \mathbf{S} (or \mathbf{R}), the diagonal elements of \mathbf{S} are equal to the diagonal elements of $\hat{\mathbf{L}}\hat{\mathbf{L}}^t + \hat{\psi}$. However, the off-diagonal elements of \mathbf{S} are not usually reproduced by $\hat{\mathbf{L}}\hat{\mathbf{L}}^t + \hat{\psi}$.

- How to determine the number of factors, m ?

Consider the residual matrix of a m factor model

$$\mathbf{S} - (\hat{\mathbf{L}}\hat{\mathbf{L}}^t + \hat{\boldsymbol{\psi}})$$

with zero diagonal elements. If the other elements are also small we will take the m factor model to be appropriate.

Ideally, the contributions of the first few factors to the sample variance should be large. The contribution to the sample variance s_{ii} from the first common factor is $\hat{\ell}_{i1}^2$. The contribution to the *total sample variance*, $s_{11} + s_{22} + \cdots + s_{pp}$, from the first common factor is

$$\sum_{i=1}^p \hat{\ell}_{i1}^2 = \left(\sqrt{\hat{\lambda}_1} \hat{\mathbf{e}}_1 \right)^t \left(\sqrt{\hat{\lambda}_1} \hat{\mathbf{e}}_1 \right) = \hat{\lambda}_1.$$

In general, the proportion of total sample variance due to the j th factor is

$$\frac{\hat{\lambda}_j}{s_{11} + s_{22} + \cdots + s_{pp}}$$

for a factor analysis of \mathbf{S} , or

$$\frac{\hat{\lambda}_j}{p}$$

for a factor analysis of \mathbf{R} .

Software packages sometimes set m equal to the number of eigenvalues of \mathbf{R} larger than 1 (if the correlation matrix is factored), or equal m to the number of positive eigenvalues of \mathbf{S} . (Be careful when using these rules blindly!)

Example: In a consumer-preference study, a number of customers were asked to rate several attributes of a new product. The correlation matrix of the responses was calculated.

Taste	1.00	0.02	0.96	0.42	0.01
Good buy for money	0.02	1.00	0.13	0.71	0.85
Flavor	0.96	0.13	1.00	0.50	0.11
Suitable for snack	0.42	0.71	0.50	1.00	0.79
Provides lots of energy	0.01	0.85	0.11	0.79	1.00

```
> library(mva)
> R <- matrix(c(1.00,0.02,0.96,0.42,0.01,
               0.02,1.00,0.13,0.71,0.85,
               0.96,0.13,1.00,0.50,0.11,
               0.42,0.71,0.50,1.00,0.79,
               0.01,0.85,0.11,0.79,1.00), 5, 5)
> eigen(R)
$values
[1] 2.85309042 1.80633245 0.20449022 0.10240947 0.03367744
```

\$vectors

	[,1]	[,2]	[,3]	[,4]	[,5]
[1,]	0.3314539	0.60721643	0.09848524	0.1386643	0.701783012
[2,]	0.4601593	-0.39003172	0.74256408	-0.2821170	0.071674637
[3,]	0.3820572	0.55650828	0.16840896	0.1170037	-0.708716714
[4,]	0.5559769	-0.07806457	-0.60158211	-0.5682357	0.001656352
[5,]	0.4725608	-0.40418799	-0.22053713	0.7513990	0.009012569

The first 2 eigenvalues of \mathbf{R} are the only ones being larger than 1. These two will account for

$$\frac{2.853 + 1.806}{5} = 0.93$$

of the total (standardized) sample variance. Thus we decide to set $m = 2$.

There is no special function available in R allowing to get the estimated factor loadings, communalities, and specific variances (uniquenesses). Hence we directly calculate those quantities.

```

> L <- matrix(rep(0, 10), 5, 2) # factor loadings
> for (j in 1:2) L[ ,j] <- sqrt(eigen(R)$values[j]) * eigen(R)$vectors[ ,j]
      [,1] [,2]
[1,] 0.560 0.816
[2,] 0.777 -0.524
[3,] 0.645 0.748
[4,] 0.939 -0.105
[5,] 0.798 -0.543
> h2 <- diag(L %*% t(L)); h2 # communalities
[1] 0.979 0.879 0.976 0.893 0.932
> psi <- diag(R) - h2; psi # specific variances
[1] 0.0205 0.1211 0.0241 0.1071 0.0678
> R - (L %*% t(L) + diag(psi)) # residuals
      [,1] [,2] [,3] [,4] [,5]
[1,] 0.0000 0.013 -0.0117 -0.020 0.0064
[2,] 0.0126 0.000 0.0205 -0.075 -0.0552
[3,] -0.0117 0.020 0.0000 -0.028 0.0012
[4,] -0.0201 -0.075 -0.0276 0.000 -0.0166
[5,] 0.0064 -0.055 0.0012 -0.017 0.0000

```


Thus we would judge a 2-factor model providing a good fit to the data. The large communalities indicate that this model accounts for a large percentage of the sample variance of each variable.

A Modified Approach – The Principle Factor Analysis

We describe this procedure in terms of a factor analysis of R . If

$$\rho = \mathbf{L}\mathbf{L}^t + \psi$$

is correctly specified, then the m common factors should account for the off-diagonal elements of ρ , as well as the communality portions of the diagonal elements

$$\rho_{ii} = 1 = h_i^2 + \psi_i.$$

If the specific factor contribution ψ_i is removed from the diagonal or, equivalently, the 1 replaced by h_i^2 the resulting matrix is $\rho - \psi = \mathbf{L}\mathbf{L}^t$.

Suppose initial estimates ψ_i^* are available. Then we replace the i th diagonal element of \mathbf{R} by $h_i^{*2} = 1 - \psi_i^*$, and obtain the reduced correlation matrix \mathbf{R}_r , which is now factored as

$$\mathbf{R}_r \approx \mathbf{L}_r^* \mathbf{L}_r^{*t} .$$

The *principle factor method* of factor analysis employs the estimates

$$\mathbf{L}_r^* = \left[\sqrt{\hat{\lambda}_1^*} \hat{\mathbf{e}}_1^*, \sqrt{\hat{\lambda}_2^*} \hat{\mathbf{e}}_2^*, \dots, \sqrt{\hat{\lambda}_m^*} \hat{\mathbf{e}}_m^* \right]$$

and

$$\hat{\psi}_i^* = 1 - \sum_{j=1}^m \ell_{ij}^{*2} ,$$

where $(\hat{\lambda}_i^*, \hat{\mathbf{e}}_i^*)$ are the (largest) eigenvalue-eigenvector pairs from \mathbf{R}_r . Re-estimate the communalities again and continue till convergence. As initial choice of h_i^{*2} you can use $1 - \psi_i^* = 1 - 1/r^{ii}$, where r^{ii} is the i th diagonal element of \mathbf{R}^{-1} .

Example cont'ed:

```
> h2 <- 1 - 1/diag(solve(R)); h2 # initial guess
[1] 0.93 0.74 0.94 0.80 0.83

> R.r <- R; diag(R.r) <- h2
> L.star <- matrix(rep(0, 10), 5, 2) # factor loadings
> for (j in 1:2) L.star[,j] <- sqrt(eigen(R.r)$values[j]) * eigen(R.r)$vectors[,j]
> h2.star <- diag(L.star %*% t(L.star)); h2.star # communalities
[1] 0.95 0.76 0.95 0.83 0.88

> # apply 3 times to get convergence

> R.r <- R; diag(R.r) <- h2.star
> L.star <- matrix(rep(0, 10), 5, 2) # factor loadings
> for (j in 1:2) L.star[,j] <- sqrt(eigen(R.r)$values[j]) * eigen(R.r)$vectors[,j]
> h2.star <- diag(L.star %*% t(L.star)); h2.star # communalities
[1] 0.97 0.77 0.96 0.83 0.93
```

```
> L.star # loadings
      [,1] [,2]
[1,] -0.60 -0.78
[2,] -0.71  0.51
[3,] -0.68 -0.71
[4,] -0.90  0.15
[5,] -0.77  0.58
```

```
> 1 - h2.star # specific variances
[1] 0.032 0.231 0.039 0.167 0.069
```

The principle components method for R can be regarded as a principal factor method with initial communality estimates of unity (or specific variance estimates equal to zero) and without iterating.

The only estimating procedure available in R is the maximum likelihood method. Beside the PCA method this is the only one, which is strongly recommended and shortly discussed now.

Maximum Likelihood Method

We now assume that the common factors F and the specific factors ϵ are from a normal distribution. Then maximum likelihood estimates of the unknown factor loadings L and the specific variances ψ may be obtained.

This strategy is the only one which is implemented in R and is now applied onto our example.

Example cont'ed:

```
> factanal(covmat = R, factors=2)
```

Call:

```
factanal(factors = 2, covmat = R, rotation = "none")
```

```
Uniquenesses: [1] 0.028 0.237 0.040 0.168 0.052
```

Loadings:

	Factor1	Factor2
[1,]	0.976	-0.139
[2,]	0.150	0.860
[3,]	0.979	
[4,]	0.535	0.738
[5,]	0.146	0.963

	Factor1	Factor2
SS loadings	2.24	2.23
Proportion Var	0.45	0.45
Cumulative Var	0.45	0.90

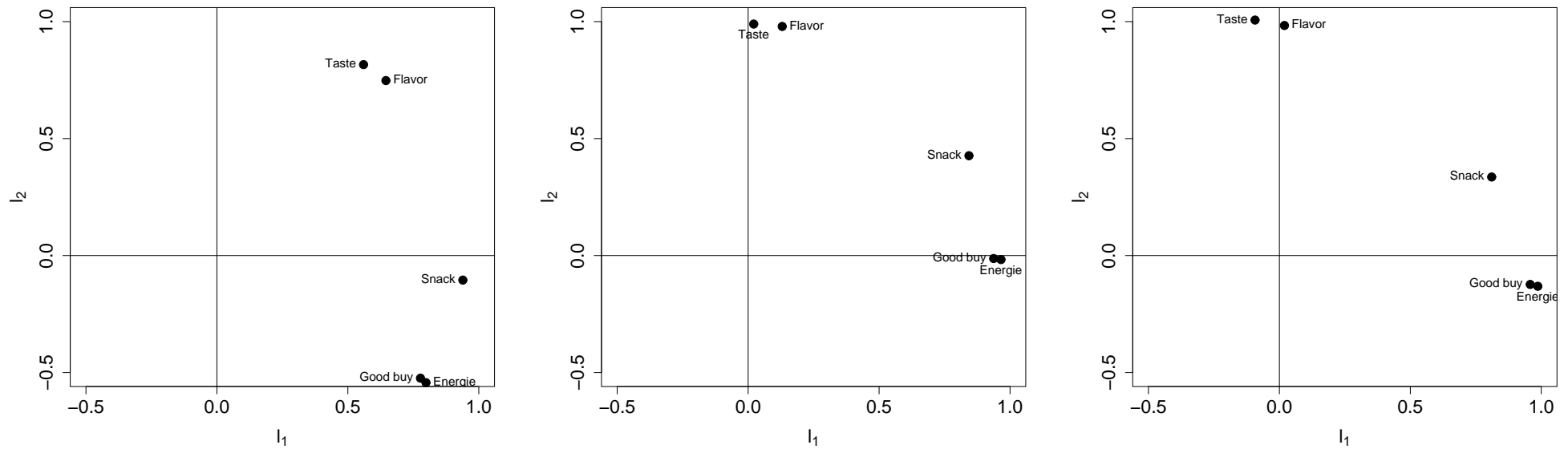
Factor Rotation

Since the original factor loadings are (a) not unique, and (b) usually not interpretable, we rotate them until a *simple structure* is achieved.

We concentrate on graphical methods for $m = 2$. A plot of the pairs of factor loadings $(\hat{\ell}_{i1}, \hat{\ell}_{i2})$, yields p points, each point corresponding to a variable. These points can be rotated by using either the *varimax* or the *promax* criterion.

Example cont'ed: Estimates of the factor loadings from the principal component approach were:

> L			> varimax(L)			> promax(L)		
	[,1]	[,2]		[,1]	[,2]		[,1]	[,2]
[1,]	0.560	0.816	[1,]	0.021	0.989	[1,]	-0.093	1.007
[2,]	0.777	-0.524	[2,]	0.937	-0.013	[2,]	0.958	-0.124
[3,]	0.645	0.748	[3,]	0.130	0.979	[3,]	0.019	0.983
[4,]	0.939	-0.105	[4,]	0.843	0.427	[4,]	0.811	0.336
[5,]	0.798	-0.543	[5,]	0.965	-0.017	[5,]	0.987	-0.131



After rotation it's much clearer to see that variables 2 (Good buy), 4 (Snack), and 5 (Energy) define factor 1 (high loadings on factor 1, small loadings on factor 2), while variables 1 (Taste) and 3 (Flavor) define factor 2 (high loadings on factor 2, small loadings on factor 1).

Johnson & Wichern call factor 1 a *nutrition factor* and factor 2 a *taste factor*.

Factor Scores

In factor analysis, interest is usually centered on the parameters in the factor model. However, the estimated values of the common factors, called *factor scores*, may also be required (e.g., for diagnostic purposes).

These scores are not estimates of unknown parameters in the usual sense. They are rather estimates of values for the unobserved random factor vectors. Two methods are provided in `factanal(..., scores =)`: the regression method of Thomson, and the weighted least squares method of Bartlett.

Both these methods allows us to plot n such p -dimensional observations as n m -dimensional scores.

Example: A factor analytic analysis of the fvc data might be as follows:

- calculate the maximum likelihood estimates of the loadings w/o rotation,
- apply a varimax rotation on these estimates and check plot of the loadings,
- estimate factor scores and plot them for the n observations.

```
> fa <- factanal(aimu[, 3:8], factors=2, scores="none", rotation="none"); fa
```

```
Uniquenesses:
```

age	height	weight	VC	FEV1	FEV1.VC
0.782	0.523	0.834	0.005	0.008	0.005

```
Loadings:
```

	Factor1	Factor2
age	-0.378	-0.274
height	0.682	-0.109
weight	0.378	-0.153
VC	0.960	-0.270
FEV1	0.951	0.295
FEV1.VC		0.993

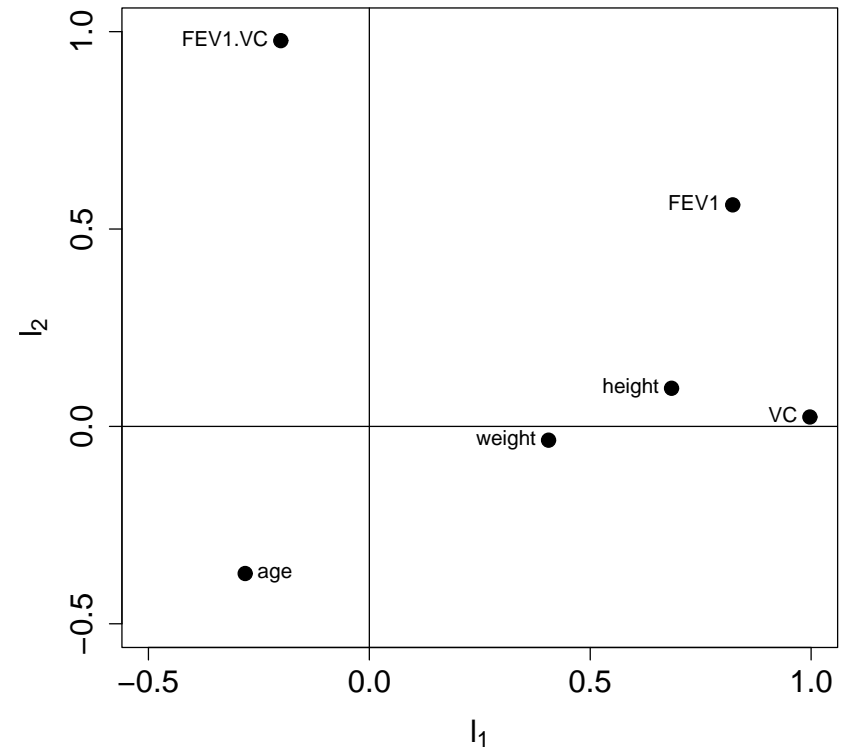
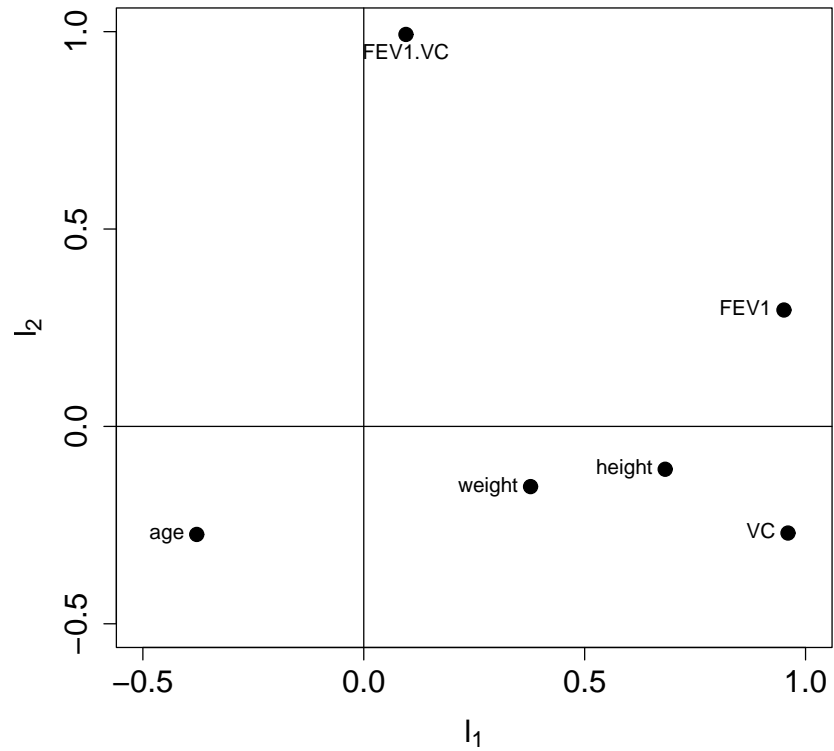
```

                Factor1 Factor2
SS loadings      2.587   1.256
Proportion Var   0.431   0.209
Cumulative Var   0.431   0.640
> L <- fa$loadings
> Lv <- varimax(fa$loadings); Lv
$loadings
      Factor1  Factor2
age      -0.2810 -0.37262
height   0.6841  0.09667
weight   0.4057 -0.03488
VC        0.9972  0.02385
FEV1      0.8225  0.56122
FEV1.VC -0.2004  0.97716

$rotmat
      [,1]  [,2]
[1,]  0.9559 0.2937
[2,] -0.2937 0.9559

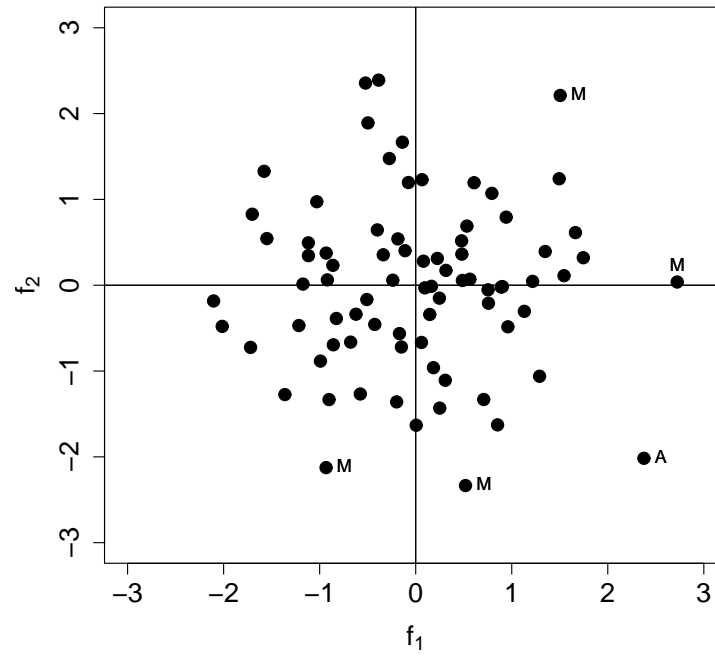
```

```
> plot(L); plot(Lv)
```



```
> s <- factanal(aimu[, 3:8], factors=2, scores="reg", rot="varimax")$scores  
> plot(s); i <- identify(s, region); aimu[i, ]
```

	nr	year	age	height	weight	VC	FEV1	FEV1.VC	region	
	25	25	85	28	189	85	740	500	68	A
	38	38	83	44	174	78	475	335	71	M
	46	46	83	23	190	75	665	635	95	M
	57	57	83	25	187	102	780	580	81	M
	71	71	83	37	173	78	590	400	68	M

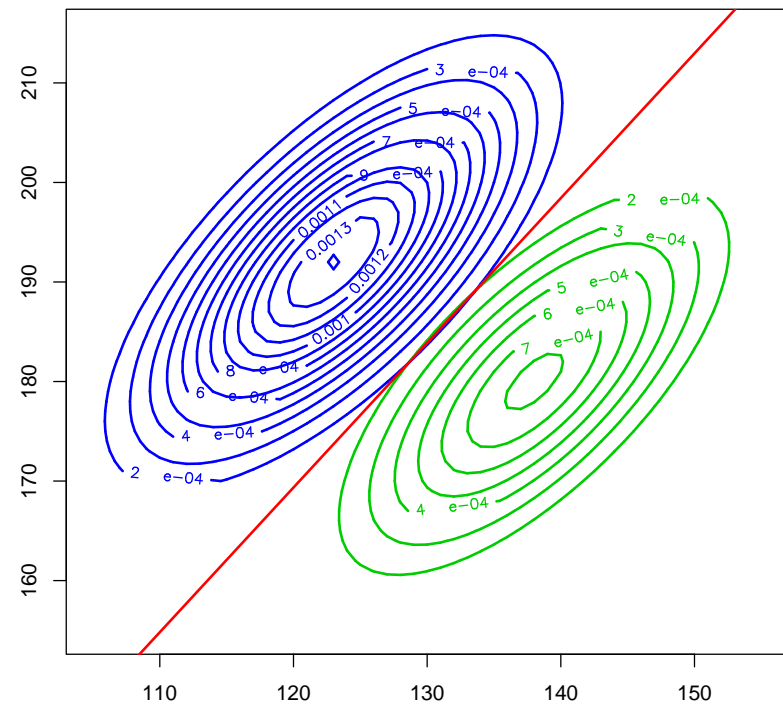
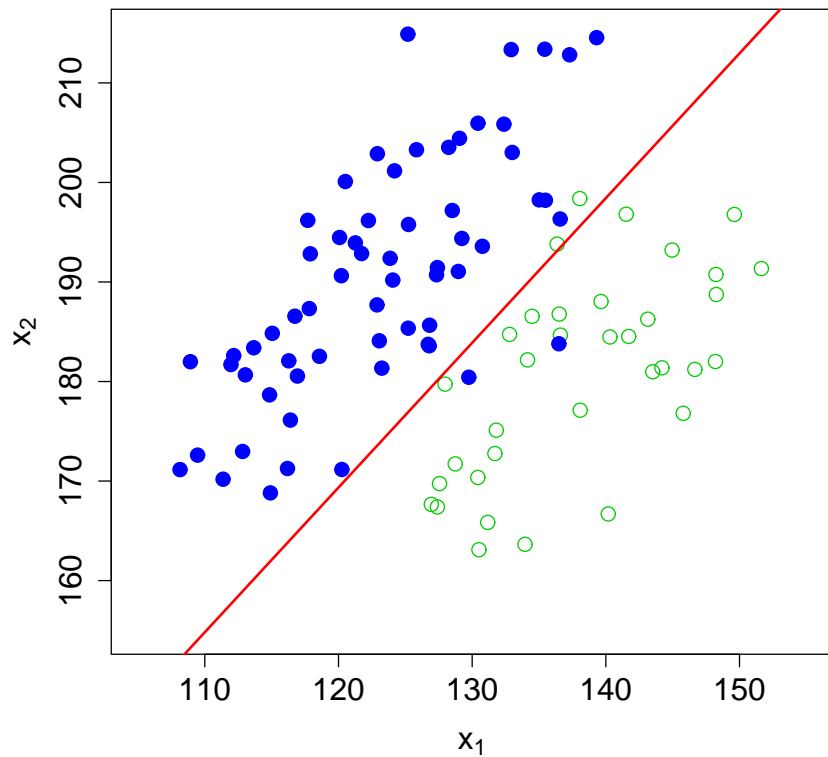


Discrimination and Classification

Discriminant analysis (DA) and classification are multivariate techniques concerned with *separating* distinct sets of objects (observations) and with *allocating* new objects to previously defined groups (defined by a categorical variable). There are several purposes for DA:

- (**Discrimination, separation**) To describe either graphically (low dimension) or algebraically, the differential features of objects from several known collections (populations, or groups).
- (**Classification, allocation**) To sort objects into 2 or more labelled classes. Thus, we derive a rule, that is used to optimally assign a *new* object to the labelled classes.

Consider 2 classes. Label these groups g_1, g_2 . The objects are to be classified on the basis of measurements on a p variate random vector $\mathbf{X} = (X_1, X_2, \dots, X_p)$. The observed values differ to some extent from one class to the other. Thus we assume that all objects \mathbf{x} in class i have density $f_i(\mathbf{x})$, $i = 1, 2$.



Technique was introduced by R.A. Fisher. His idea was to transform the multivariate \mathbf{x} to univariate y such that the y 's derived from population g_1 and g_2 were separated as much as possible. He considered linear combinations of \mathbf{x} .

If we let μ_{1Y} be the mean of Y obtained from \mathbf{X} belonging to g_1 , and μ_{2Y} be the mean of Y obtained from \mathbf{X} belonging to g_2 , then he selected the linear combination that maximized the squared distance between μ_{1Y} and μ_{2Y} relative to the variability of the Y 's.

We define

$$\boldsymbol{\mu}_1 = \mathbf{E}(\mathbf{X}|g_1), \quad \text{and} \quad \boldsymbol{\mu}_2 = \mathbf{E}(\mathbf{X}|g_2)$$

and suppose the covariance matrix

$$\boldsymbol{\Sigma} = \mathbf{E}((\mathbf{X} - \boldsymbol{\mu}_i)(\mathbf{X} - \boldsymbol{\mu}_i)^t), \quad i = 1, 2$$

is the same for both populations (somewhat critical in practice).

We consider the linear combination

$$Y = \boldsymbol{\ell}^t \mathbf{X}$$

and get population-specific means

$$\mu_{1Y} = \mathbb{E}(Y|g_1) = \mathbb{E}(\boldsymbol{\ell}^t \mathbf{X}|g_1) = \boldsymbol{\ell}^t \boldsymbol{\mu}_1$$

$$\mu_{2Y} = \mathbb{E}(Y|g_2) = \mathbb{E}(\boldsymbol{\ell}^t \mathbf{X}|g_2) = \boldsymbol{\ell}^t \boldsymbol{\mu}_2$$

but equal variance

$$\sigma_Y^2 = \text{var}(Y) = \text{var}(\boldsymbol{\ell}^t \mathbf{X}) = \boldsymbol{\ell}^t \text{cov}(\mathbf{X}) \boldsymbol{\ell} = \boldsymbol{\ell}^t \boldsymbol{\Sigma} \boldsymbol{\ell}.$$

The best linear combination is derived from the ratio ($\boldsymbol{\delta} = \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$)

$$\frac{(\mu_{1Y} - \mu_{2Y})^2}{\sigma_Y^2} = \frac{(\boldsymbol{\ell}^t \boldsymbol{\mu}_1 - \boldsymbol{\ell}^t \boldsymbol{\mu}_2)^2}{\boldsymbol{\ell}^t \boldsymbol{\Sigma} \boldsymbol{\ell}} = \frac{\boldsymbol{\ell}^t (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \boldsymbol{\ell}}{\boldsymbol{\ell}^t \boldsymbol{\Sigma} \boldsymbol{\ell}} = \frac{(\boldsymbol{\ell}^t \boldsymbol{\delta})^2}{\boldsymbol{\ell}^t \boldsymbol{\Sigma} \boldsymbol{\ell}}$$

Result: Let $\delta = \mu_1 - \mu_2$ and $Y = \ell^t \mathbf{X}$, then

$$\frac{(\ell^t \delta)^2}{\ell^t \Sigma \ell}$$

is maximized by the choice

$$\ell = c \Sigma^{-1} \delta = c \Sigma^{-1} (\mu_1 - \mu_2)$$

for any $c \neq 0$. Choosing $c = 1$ produces the linear combination

$$Y = \ell^t \mathbf{X} = (\mu_1 - \mu_2)^t \Sigma^{-1} \mathbf{X}$$

which is known as *Fisher's linear discriminant function*.

We can also employ this result as *classification device*. Let $y_0 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \boldsymbol{\Sigma}^{-1} \boldsymbol{x}_0$ be the value of the discriminant function for a new observation \boldsymbol{x}_0 and let

$$m = \frac{1}{2}(\mu_{1Y} + \mu_{2Y}) = \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)$$

be the midpoint of the 2 univariate population means. It can be shown that

$$E(Y_0|g_1) - m \geq 0 \quad \text{and} \quad E(Y_0|g_2) - m < 0$$

That is, if \boldsymbol{X}_0 is from g_1 , Y_0 is expected to be larger than the midpoint. If \boldsymbol{X}_0 is from g_2 , Y_0 is expected to be smaller. Thus the classification rule is:

$$\text{Allocate } \boldsymbol{x}_0 \text{ to } g_1 \text{ if: } \quad y_0 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \boldsymbol{\Sigma}^{-1} \boldsymbol{x}_0 \geq m$$

$$\text{Allocate } \boldsymbol{x}_0 \text{ to } g_2 \text{ if: } \quad y_0 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \boldsymbol{\Sigma}^{-1} \boldsymbol{x}_0 < m .$$

Because the population moments are not known, we replace μ_1 , μ_2 , and Σ by their empirical versions.

Suppose we have 2 data matrices \mathbf{X}_1 from g_1 and \mathbf{X}_2 from g_2 with n_1 and n_2 observations, from which we calculate both sample means $\bar{\mathbf{x}}_1$, $\bar{\mathbf{x}}_2$, and sample covariance matrices \mathbf{S}_1 , \mathbf{S}_2 . Since it is assumed that the covariance matrices in the groups are the same, we combine (pool) \mathbf{S}_1 and \mathbf{S}_2 to derive a single estimate of Σ . Hence we use the *pooled sample covariance* matrix

$$\mathbf{S}_p = \frac{(n_1 - 1)\mathbf{S}_1 + (n_2 - 1)\mathbf{S}_2}{n_1 + n_2 - 2}$$

an unbiased estimate of Σ . Now, μ_1 , μ_2 , and Σ are replaced by $\bar{\mathbf{x}}_1$, $\bar{\mathbf{x}}_2$, and \mathbf{S}_p in the previous formulas to give *Fisher's sample linear discriminant function*

$$y = \hat{\ell}^t \mathbf{x} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^t \mathbf{S}_p^{-1} \mathbf{x}.$$

The midpoint between both sample means is

$$\hat{m} = \frac{1}{2}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^t \mathbf{S}_p^{-1}(\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2)$$

and the classification rule becomes

$$\text{Allocate } \mathbf{x}_0 \text{ to } g_1 \text{ if: } (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^t \mathbf{S}_p^{-1} \mathbf{x}_0 \geq \hat{m}$$

$$\text{Allocate } \mathbf{x}_0 \text{ to } g_2 \text{ if: } (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^t \mathbf{S}_p^{-1} \mathbf{x}_0 < \hat{m}.$$

This idea can be easily generalized onto more than 2 classes. Moreover, instead of using a linear discriminant function we can also use a quadratic one.

Example: Fisher's Iris data

Data describing the sepal (Kelchblatt) width and length, and the petal (Blütenblatt) width and length of 3 different Iris species (Setosa, Versicolor, Virginica) were observed. There are 50 observations for each species.

```
> library(MASS)
> data(iris3)
> Iris <- data.frame(rbind(iris3[, ,1], iris3[, ,2], iris3[, ,3]),
                    Sp = rep(c("s","c","v"), rep(50,3)))
> z <- lda(Sp ~ Sepal.L.+Sepal.W.+Petal.L.+Petal.W., Iris, prior = c(1,1,1)/3)
Prior probabilities of groups:
              c              s              v
0.33333333 0.33333333 0.33333333

Group means:
  Sepal.L. Sepal.W. Petal.L. Petal.W.
c    5.936    2.770    4.260    1.326
s    5.006    3.428    1.462    0.246
v    6.588    2.974    5.552    2.026
```

Coefficients of linear discriminants:

	LD1	LD2
Sepal.L.	-0.8293776	0.02410215
Sepal.W.	-1.5344731	2.16452123
Petal.L.	2.2012117	-0.93192121
Petal.W.	2.8104603	2.83918785

Proportion of trace:

LD1	LD2
0.9912	0.0088

> predict(z, Iris)\$class

```
[1] s s s s s s s s s s s s s s s s s s s s s s s s s s s s
[32] s s s s s s s s s s s s s s s s s s c c c c c c c c c c
[63] c c c c c c c c v c c c c c c c c c c c c v c c c c c c
[94] c c c c c c c v v v v v v v v v v v v v v v v v v v v v
[125] v v v v v v v v v c v v v v v v v v v v v v v v v v
```

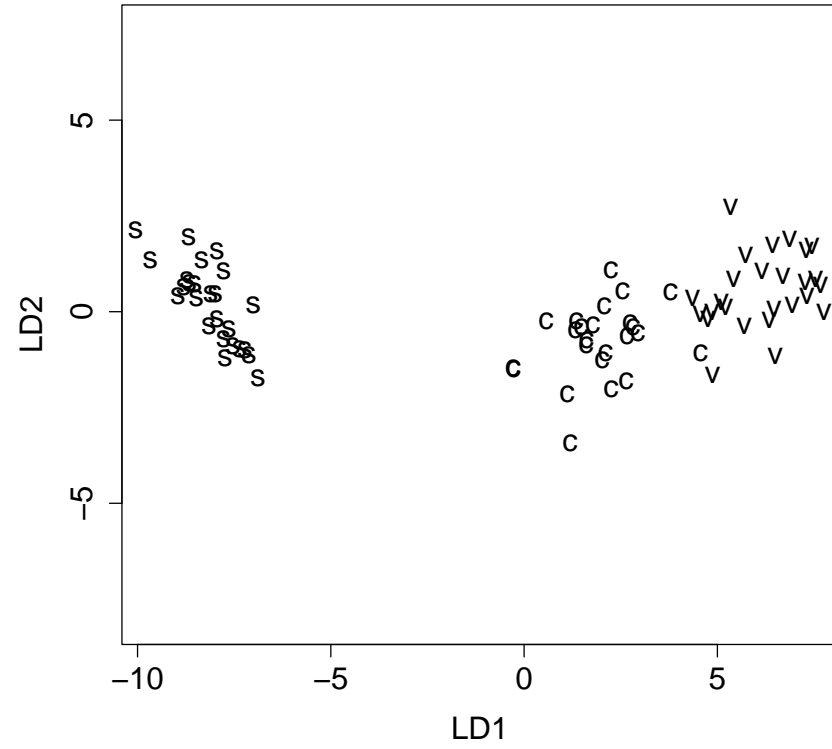
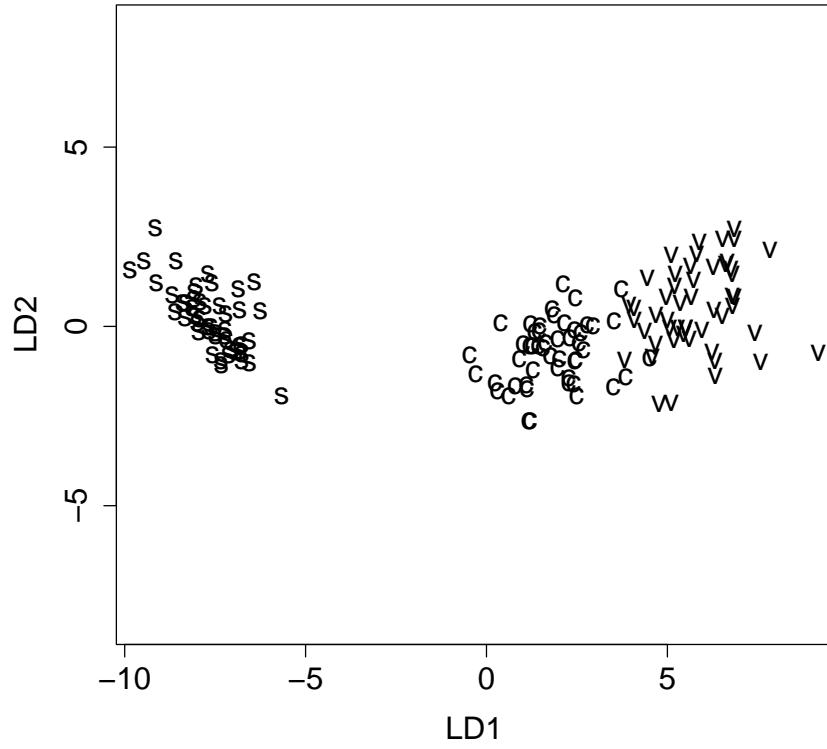
```

> table(predict(z, Iris)$class, Iris$Sp)
      c  s  v
c 48  0  1
s  0 50  0
v  2  0 49
> train <- sample(1:150, 75); table(Iris$Sp[train])
      c  s  v
24 25 26
> z1 <- lda(Sp ~ Sepal.L.+Sepal.W.+Petal.L.+Petal.W., Iris,
            prior = c(1,1,1)/3, subset = train)
> predict(z1, Iris[-train, ])$class
 [1] s s s s s s s s s s s s s s s s s s s s s s s s s c c c c c c
[32] c c c c c c v c c c c c c c c c c c c c v v v v v v v v v v v
[63] v v v c v v v v v v v v v v
> table(predict(z1, Iris[-train, ])$class, Iris[-train, ]$Sp)
      c  s  v
c 25  0  1
s  0 25  0
v  1  0 23

```



```
> plot(z); plot(z1)
```



```

> ir.ld <- predict(z, Iris)$x # => LD1 and LD2 coordinates
> eqscplot(ir.ld, type="n", xlab="First LD", ylab="Second LD") # eq. scaled axes
> text(ir.ld, as.character(Iris$Sp)) # plot LD1 vs. LD2

> # calc group-spec. means of LD1 & LD2
> tapply(ir.ld[ , 1], Iris$Sp, mean)
      c      s      v
1.825049 -7.607600  5.782550
> tapply(ir.ld[ , 2], Iris$Sp, mean)
      c      s      v
-0.7278996  0.2151330  0.5127666

> # faster alternative:
> ir.m <- lda(ir.ld, Iris$Sp)$means; ir.m
      LD1      LD2
c  1.825049 -0.7278996
s -7.607600  0.2151330
v  5.782550  0.5127666
> points(ir.m, pch=3, mkh=0.3, col=2) # plot group means as "+"

```

```

> perp <- function(x, y, ...) {
+   m <- (x+y)/2           # midpoint of the 2 group means
+   s <- -(x[1]-y[1])/(x[2]-y[2]) # perpendicular line through midpoint
+   abline(c(m[2]-s*m[1], s), ...) # draw classification regions
+   invisible()
> }
> perp(ir.m[1,], ir.m[2,], col=1) # classification decision b/w groups 1&2
> perp(ir.m[1,], ir.m[3,], col=2) # classification decision b/w groups 1&3
> perp(ir.m[2,], ir.m[3,], col=3) # classification decision b/w groups 2&3

```

