

# Honours Multivariate Analysis

Lecture 1 - Introduction, Examples of Multivariate Data, Visualization and Summary Statistics

Stefan S. Britz

Department of Statistical Sciences  
University of Cape Town



# Course Outline

- ① Introduction, Examples of Multivariate Data †
- ② TOOLS
  - ① Visualization and Summary Statistics †
  - ② Singular Value Decomposition, Eigenvalue Decomposition and Spectral Decomposition revisited †
  - ③ The Multivariate Normal Distribution †
  - ④ Multivariate Maximum Likelihood Estimation †
  - ⑤ Multivariate Inference †
- ③ EXPLORATORY ANALYSIS
  - ① Principal Component Analysis §
  - ② Factor Analysis §
  - ③ Correspondence Analysis §
- ④ CONFIRMATORY ANALYSIS
  - ① For grouped Multivariate Data:
    - ① Manova †
    - ② Discriminant Analysis §
  - ② Regression
    - ① Multivariate Regression §
    - ② Canonical Correlation Analysis §

†Mr Stefan Britz

§Mr Miguel Rodo

# Introduction



Multivariate Analysis is concerned with the analysis of multiple response variables.  
For example,

- ① Mass (in grams), snout-vent length (in mm) and hind limb span (in mm) measurements for lizards.
- ② Air pollution data consisting of 42 measurements of 7 air-pollution variables recorded on different days.
- ③ National Track records for men and women.

# Examples of datasets

```
head(lizards)
```

	Mass	SVL	HLS
1	5.526	59.0	113.5
2	10.401	75.0	142.0
3	9.213	69.0	124.0
4	8.953	67.5	125.0
5	7.063	62.0	129.5
6	6.610	62.0	123.0

```
head(AirPollution)
```

	Wind	Radiation	CO	NO	NO2	O3	HC
1	8	98	7	2	12	8	2
2	7	107	4	3	9	5	3
3	7	103	4	3	5	6	3
4	10	88	5	2	8	15	4
5	6	91	4	2	8	10	3
6	8	90	5	2	12	12	4

# Examples of datasets

```
head(Track)
```

	Country	m100	m200	m400	m800	m1500	m3000	Marathon	Gender
1	ARG	11.57	22.94	52.50	2.05	4.25	9.19	150.32	Female
2	AUS	11.12	22.23	48.63	1.98	4.02	8.63	143.51	Female
3	AUT	11.15	22.70	50.62	1.94	4.05	8.78	154.35	Female
4	BEL	11.14	22.48	51.45	1.97	4.08	8.82	143.05	Female
5	BER	11.46	23.05	53.30	2.07	4.29	9.81	174.18	Female
6	BRA	11.17	22.60	50.62	1.97	4.17	9.04	147.41	Female

# Matrix representation of data

Data can be represented in the form of a matrix where the rows reflect the observations and the columns reflect the variables. Typically

$n = \text{Number of observations}$

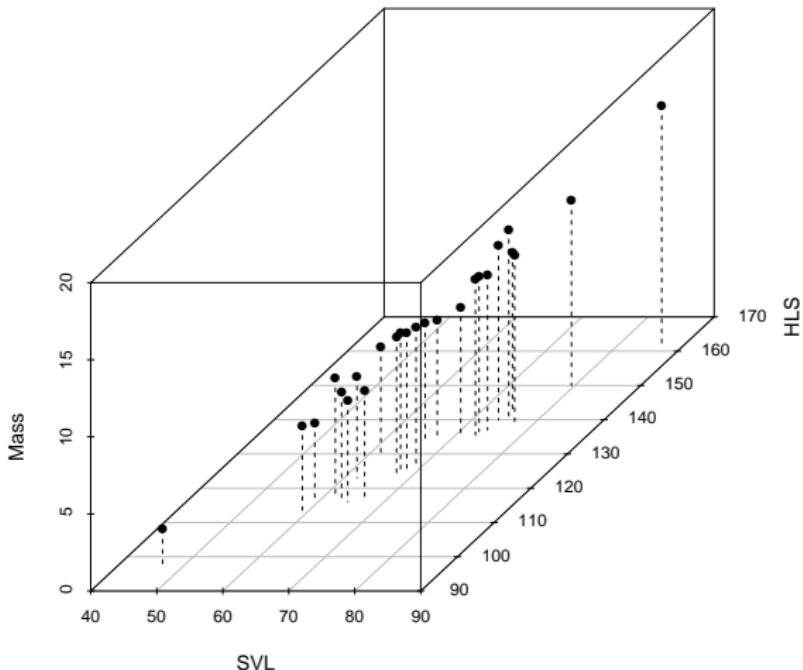
$p = \text{Number of variables}$

$$\mathbf{X}_{n \times p} = \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1p} \\ X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & & \vdots \\ X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix}$$

So for the lizard data,

$$\mathbf{X}_{25 \times 3} = \begin{bmatrix} 5.526 & 59.0 & 113.5 \\ 10.401 & 75.0 & 142.0 \\ \vdots & \vdots & \vdots \\ 6.890 & 63.0 & 117.0 \end{bmatrix}$$

The data can be viewed as  $n$  points in  $p$ -dimensional space



# Summary Statistics: Mean Vector

Similar to the univariate case, we can define a sample **mean** and variance, as well as covariances and correlations.

$$\bar{\mathbf{x}}_{p \times 1} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix}$$

where  $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$  is the sample mean of the  $j^{th}$  variable ( $j^{th}$  column of  $\mathbf{X}$ )

Can be written as

$$\bar{\mathbf{x}} = \frac{1}{n} \mathbf{X}' \mathbf{1}$$

where  $\mathbf{1}_{n \times 1}$  is a vector of ones

# Summary Statistics: Sample Mean Vector

**Example:**

If

$$\mathbf{X}_{3 \times 2} = \begin{bmatrix} 4 & 1 \\ -1 & 3 \\ 3 & 5 \end{bmatrix}$$

then

$$\frac{1}{n} \times \mathbf{X}' \mathbf{1} = \frac{1}{3} \begin{bmatrix} 4 & -1 & 3 \\ 1 & 3 & 5 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 6/3 = 2 \\ 9/3 = 3 \end{bmatrix}$$

# Summary Statistics: Sample Covariance Matrix

Similar to the univariate case, we can define a sample mean and **variance**, as well as **covariances** and correlations.

$S_{p \times p}$  is a symmetric matrix with

$$s_{jj} = \text{Var}(X_j) = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$$

$$s_{jk} = \text{Cov}(X_j, X_k) = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$$

## Summary Statistics: Sample Covariance Matrix

In matrix notation, the **sample covariance matrix** can be expressed in terms of:

$$\mathbf{1}\bar{x}' = \frac{1}{n} \mathbf{1}\mathbf{1}' \mathbf{X} \rightarrow \text{an } n \times p \text{ matrix of means, and}$$

$$\mathbf{X} - \frac{1}{n} \mathbf{1}\mathbf{1}' \mathbf{X} \rightarrow \text{an } n \times p \text{ matrix of deviations such that:}$$

# Summary Statistics: Sample Covariance Matrix

In matrix notation, the **sample covariance matrix** can be expressed in terms of:

$$\mathbf{1}\bar{\mathbf{x}}' = \frac{1}{n} \mathbf{1}\mathbf{1}' \mathbf{X} \rightarrow \text{an } n \times p \text{ matrix of means, and}$$

$$\mathbf{X} - \frac{1}{n} \mathbf{1}\mathbf{1}' \mathbf{X} \rightarrow \text{an } n \times p \text{ matrix of deviations such that:}$$

$$\mathbf{S} = \frac{1}{n-1} \left[ \left( \mathbf{X} - \frac{1}{n} \mathbf{1}\mathbf{1}' \mathbf{X} \right)' \left( \mathbf{X} - \frac{1}{n} \mathbf{1}\mathbf{1}' \mathbf{X} \right) \right]$$

$$\begin{aligned}(n-1)\mathbf{S} &= \left( \mathbf{X} - \frac{1}{n} \mathbf{1}\mathbf{1}' \mathbf{X} \right)' \left( \mathbf{X} - \frac{1}{n} \mathbf{1}\mathbf{1}' \mathbf{X} \right) \\&= \mathbf{X}' \left( \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}' \right)' \left( \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}' \right) \mathbf{X} \\&= \mathbf{X}' \left( \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}' \right) \mathbf{X} \\&= \mathbf{X}' \mathbf{X} - \frac{1}{n} \mathbf{X}' \mathbf{1}\mathbf{1}' \mathbf{X} \\&= \mathbf{X}' \mathbf{X} - n \bar{\mathbf{x}} \bar{\mathbf{x}}'\end{aligned}$$

## Summary Statistics: Sample Correlation Matrix

Similar to the univariate case, we can define a sample mean and variance, as well as covariances and **correlations**.

$$\mathbf{R}_{p \times p} = \begin{bmatrix} 1 & r_{12} & \cdots & r_{1p} \\ r_{21} & 1 & \cdots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ r_{p1} & r_{p2} & \cdots & 1 \end{bmatrix}$$

where

$$r_{jk} = \frac{s_{jk}}{\sqrt{s_{jj}s_{kk}}}$$

Can be written as

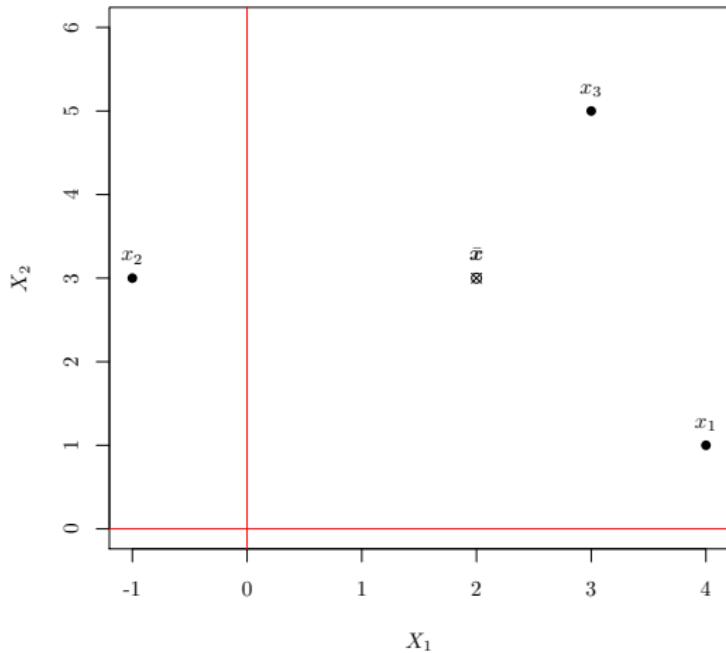
$$\mathbf{R} = [\text{diag}(\mathbf{S})]^{-\frac{1}{2}} \mathbf{S} [\text{diag}(\mathbf{S})]^{-\frac{1}{2}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{S} \mathbf{D}^{-\frac{1}{2}}$$

# Summary Statistics: Applications in R

See [Lecture1\\_1.R](#)

# Graphical interpretation of a matrix

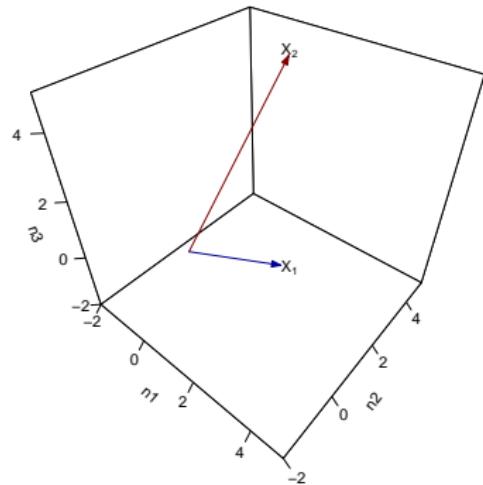
An  $n \times p$  matrix  $\mathbf{X}$  can be viewed as  $n$  observations in  $p$ -dimensional space...



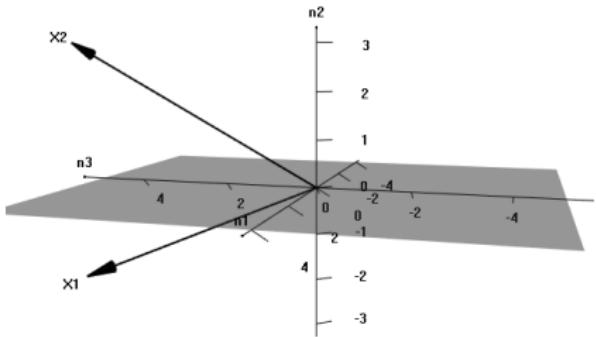
# Graphical interpretation of a matrix

...or as  $p$  vectors in  $n$ -dimensional space.

Using `plot3D::arrows3D`



Using `matplotlib::arrows3d & rgl`



# Deviation Vectors

We note that the expressions for the covariance matrix:

$$S = \frac{1}{n-1} \left[ \left( \mathbf{X} - \frac{1}{n} \mathbf{1} \mathbf{1}' \mathbf{X} \right)' \left( \mathbf{X} - \frac{1}{n} \mathbf{1} \mathbf{1}' \mathbf{X} \right) \right]$$

depend on the deviation of the observed values from the variable means. We can define and calculate so-called "deviation vectors" and then calculate the elements of the covariance matrix using these deviation vectors. The graphical display and interpretation of the deviation vectors is also informative.

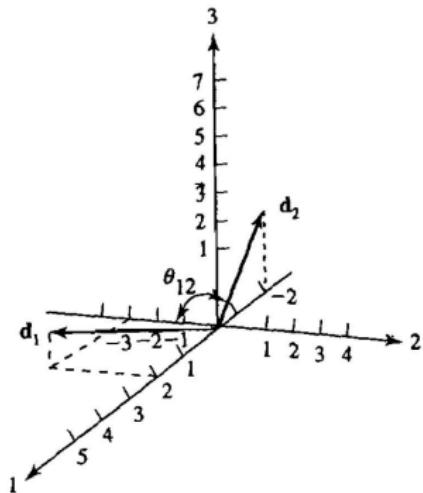
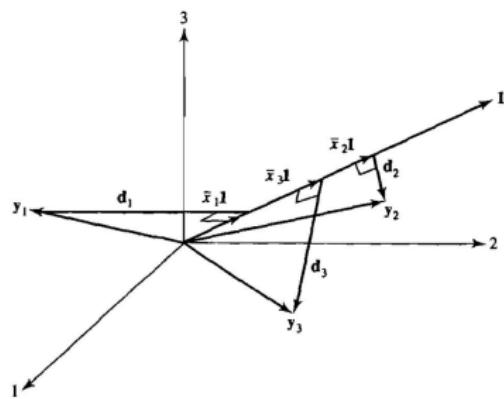
For an illustration, continuing with the previous example, see again

Lecture1\_1.R .

## Geometry of deviation vectors

- ① The projection of a column  $x_i$  of the data matrix  $\mathbf{X}$  into the equal angular vector  $\mathbf{1}$  is the vector  $\bar{x}_i \mathbf{1}$  with length  $\sqrt{n}|\bar{x}_i|$ , i.e., the  $i^{\text{th}}$  sample mean is related to the projection of  $x_i$  on  $\mathbf{1}$ .
- ② The information comprising  $\mathbf{S}$  is obtained from the deviation vectors  $d_i = x_i - \bar{x}_i \mathbf{1}$ . The square of the length of  $d_i$  is  $(n_i - 1)s_{ii}$  and the inner product between  $d_i$  and  $d_k$  is  $(n - 1)s_{ik}$ .
- ③ The sample correlation coefficient  $r_{ik}$  is the cosine of the angle between  $d_i$  and  $d_k$ .

# Geometry of deviation vectors



# Homework exercise 1

Johnson & Wichern exercise 3.1.

Given the data matrix

$$\mathbf{X} = \begin{bmatrix} 9 & 1 \\ 5 & 3 \\ 1 & 2 \end{bmatrix}$$

- ① Graph the scatterplot in  $p = 2$  dimensions. Locate the sample mean on your diagram.
- ② Sketch the  $n = 3$  dimensional representation of the data, and plot the two deviation vectors.
- ③ Sketch the deviation vectors in (2) emanating from the origin. Calculate the lengths of these vectors and the cosine of the angle between them. Relate these quantities to  $\mathbf{S}$  and  $\mathbf{R}$ .

# Random Vectors

Let  $\mathbf{X}_{p \times 1}$  be a random vector whose components are random variables

$$\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix}$$

where each random variable  $X_j$  has its own marginal probability distribution  $f_{X_j}(x)$ .

# Expected value of $X$

$$\boldsymbol{\mu}_{p \times 1} = E(\mathbf{X}) = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_p) \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix}$$

where

$$\mu_j = E(X_j) = \begin{cases} \int_{-\infty}^{\infty} x f_{X_j}(x) dx \\ \sum_x x \Pr(X_j = x) \end{cases}$$

# Covariance & Correlation Matrix of $\mathbf{X}$

The covariance matrix  $\Sigma_{p \times p}$  is symmetric, positive semi-definite, and contains the pairwise covariances:

$$\Sigma = \text{Var}(\mathbf{X}) = E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})']$$

with

$$\sigma_{jj} = \text{Var}(X_j) = E[(X_j - \mu_j)^2]$$

and

$$\sigma_{jk} = \text{Cov}(X_j, X_k) = E[(X_j - \mu_j)(X_k - \mu_k)]$$

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \cdots & \sigma_{pp} \end{bmatrix}$$

# Covariance & Correlation Matrix of $\mathbf{X}$

Similarly, the correlation matrix  $\mathbf{P}_{p \times p}$  contains the pairwise correlations:

$$\mathbf{P} = \begin{bmatrix} 1 & \rho_{12} & \cdots & \rho_{1p} \\ \rho_{21} & 1 & \cdots & \rho_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{p1} & \rho_{p2} & \cdots & 1 \end{bmatrix}$$

which can be written as

$$\mathbf{P} = [diag(\Sigma)]^{-\frac{1}{2}} \Sigma [diag(\Sigma)]^{-\frac{1}{2}}$$

# Moments of Sample Mean and Covariance Matrix

We consider the rows of the  $\mathbf{X}$  matrix to represent a random sample from the joint distribution of the  $p$  variables with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ .

Therefore,  $\bar{\mathbf{X}}$  and  $\mathbf{S}$  are random variables! We'll look at their sampling distributions in detail later, for now we'll focus on their first central moments:

$$\mathrm{E}(\bar{\mathbf{X}}) = \boldsymbol{\mu} \text{ and } \mathrm{Cov}(\bar{\mathbf{X}}) = \frac{1}{n} \boldsymbol{\Sigma}$$

$$\mathrm{E}(\mathbf{S}) = \boldsymbol{\Sigma}$$

Therefore,

$$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})'$$

is an unbiased estimator of  $\boldsymbol{\Sigma}$ .

# Linear Combinations of Variables

Let  $\mathbf{X}_{p \times 1}$  be a random vector with  $E(\mathbf{X}) = \boldsymbol{\mu}$  and  $\text{Var}(\mathbf{X}) = \boldsymbol{\Sigma}$ . Now consider the following two linear combinations of the variables contained in  $\mathbf{X}$ :

$$\mathbf{b}'\mathbf{X} = b_1X_1 + b_2X_2 + \dots + b_pX_p \text{ and } \mathbf{c}'\mathbf{X} = c_1X_1 + c_2X_2 + \dots + c_pX_p.$$

Then

- ① the mean of  $\mathbf{b}'\mathbf{X}$  is  $E(\mathbf{b}'\mathbf{X}) = \mathbf{b}'\boldsymbol{\mu}$ . Likewise  $E(\mathbf{c}'\mathbf{X}) = \mathbf{c}'\boldsymbol{\mu}$ .
- ② the variance of  $\mathbf{b}'\mathbf{X}$  is  $\text{Var}(\mathbf{b}'\mathbf{X}) = \mathbf{b}'\boldsymbol{\Sigma}\mathbf{b}$ . Likewise  $\text{Var}(\mathbf{c}'\mathbf{X}) = \mathbf{c}'\boldsymbol{\Sigma}\mathbf{c}$ .
- ③ the covariance of  $\mathbf{b}'\mathbf{X}$  and  $\mathbf{c}'\mathbf{X}$  is  $\text{Cov}(\mathbf{b}'\mathbf{X}, \mathbf{c}'\mathbf{X}) = \mathbf{b}'\boldsymbol{\Sigma}\mathbf{c}$ .

# Linear Combinations of Variables

Let  $\mathbf{X}_{p \times 1}$  be a random vector with  $E(\mathbf{X}) = \boldsymbol{\mu}$  and  $\text{Var}(\mathbf{X}) = \boldsymbol{\Sigma}$ . Now consider the following two linear combinations of the variables contained in  $\mathbf{X}$ :  
 $\mathbf{b}'\mathbf{X} = b_1X_1 + b_2X_2 + \dots + b_pX_p$  and  $\mathbf{c}'\mathbf{X} = c_1X_1 + c_2X_2 + \dots + c_pX_p$ .

Then

- ① the mean of  $\mathbf{b}'\mathbf{X}$  is  $E(\mathbf{b}'\mathbf{X}) = \mathbf{b}'\boldsymbol{\mu}$ . Likewise  $E(\mathbf{c}'\mathbf{X}) = \mathbf{c}'\boldsymbol{\mu}$ .
- ② the variance of  $\mathbf{b}'\mathbf{X}$  is  $\text{Var}(\mathbf{b}'\mathbf{X}) = \mathbf{b}'\boldsymbol{\Sigma}\mathbf{b}$ . Likewise  $\text{Var}(\mathbf{c}'\mathbf{X}) = \mathbf{c}'\boldsymbol{\Sigma}\mathbf{c}$ .
- ③ the covariance of  $\mathbf{b}'\mathbf{X}$  and  $\mathbf{c}'\mathbf{X}$  is  $\text{Cov}(\mathbf{b}'\mathbf{X}, \mathbf{c}'\mathbf{X}) = \mathbf{b}'\boldsymbol{\Sigma}\mathbf{c}$ .

Homework exercise: Johnson & Wichern exercise 3.18

Energy consumption by state from the major sources,  $x_1$  = petroleum,  $x_2$  = natural gas,  $x_3$  = hydroelectric power and  $x_4$  = nuclear electric power is recorded in quadrillions ( $10^{15}$ ).

## Homework exercise 2

The resulting mean and covariance matrix are

$$\bar{x} = \begin{bmatrix} 0.766 \\ 0.508 \\ 0.438 \\ 0.161 \end{bmatrix}$$

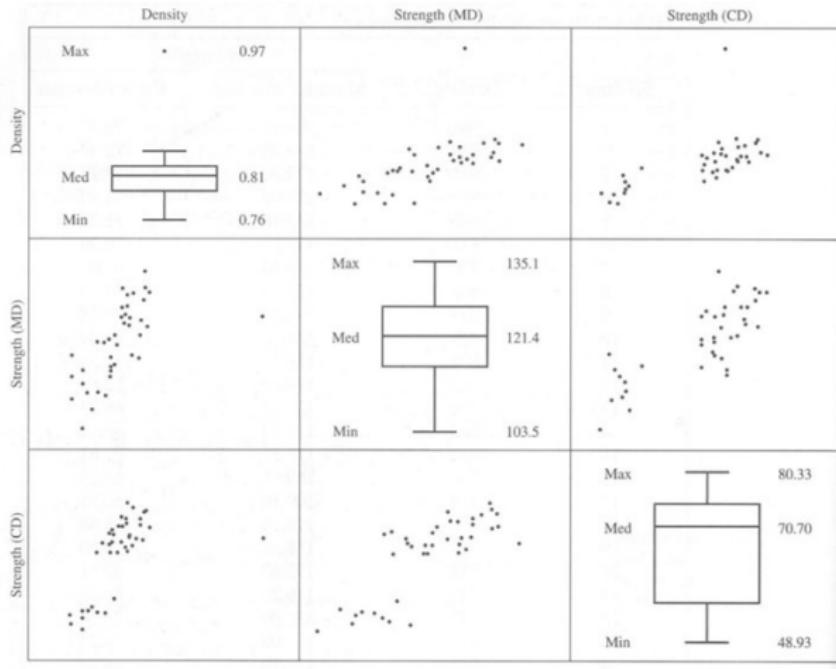
and

$$S = \begin{bmatrix} 0.856 & 0.635 & 0.173 & 0.096 \\ 0.635 & 0.568 & 0.128 & 0.067 \\ 0.173 & 0.128 & 0.171 & 0.039 \\ 0.096 & 0.067 & 0.039 & 0.043 \end{bmatrix}$$

- ① Determine the sample mean and variance of a state's total energy consumption for these major sources.
- ② Determine the sample mean and variance of the excess of petroleum consumption over natural gas consumption. Also find the sample covariance of this variable with the total variable in part 1.

# Graphical Summaries and Exploratory Plots in R

Function: `pairs(...)` #See `Lecture1_2.R` for examples on usage



# Graphical Displays in R

Package: rgl

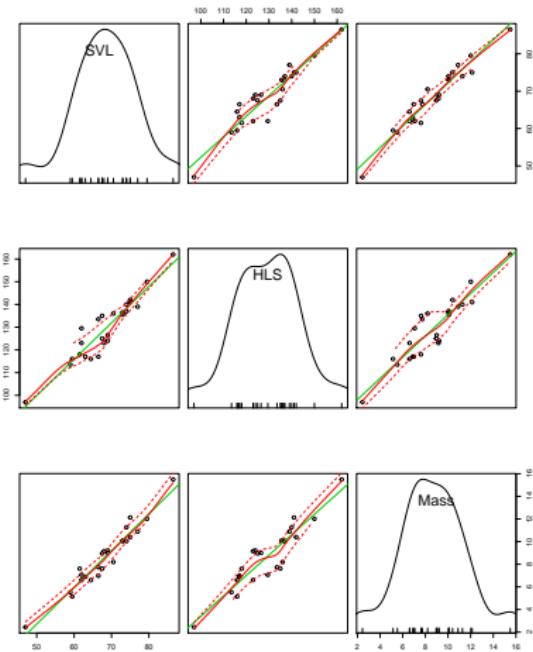
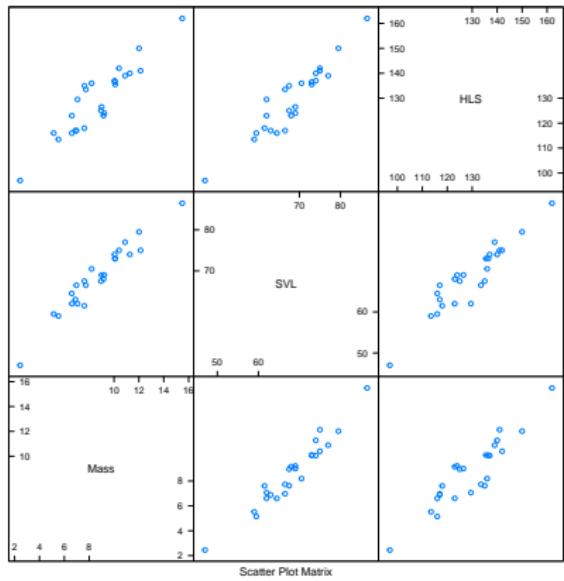
Functions:

- plot3d
- spheres3d
- points3d
- planes3d

Note that there are  
many other R  
libraries to produce  
plots, this is just my  
favourite!

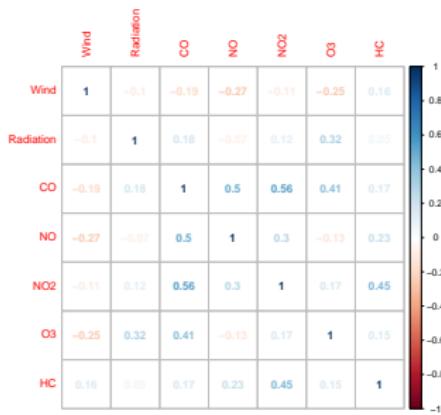
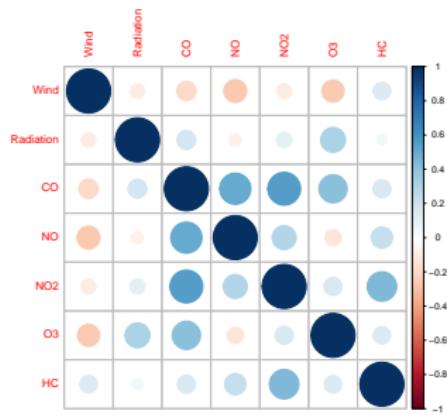
# Exploratory plots for multivariate data

See [Lecture1\\_2.R](#) for scatterplot examples.



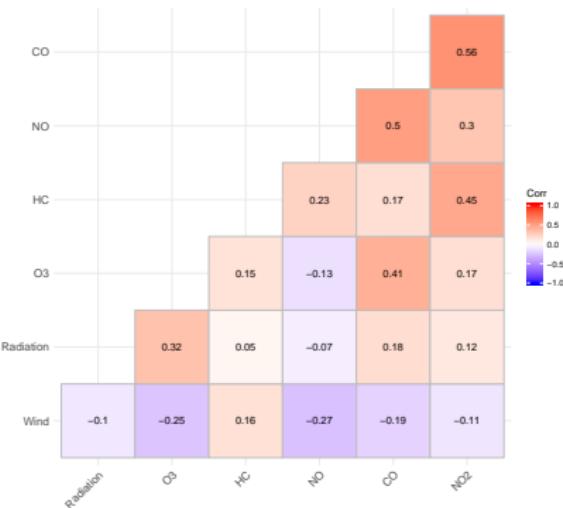
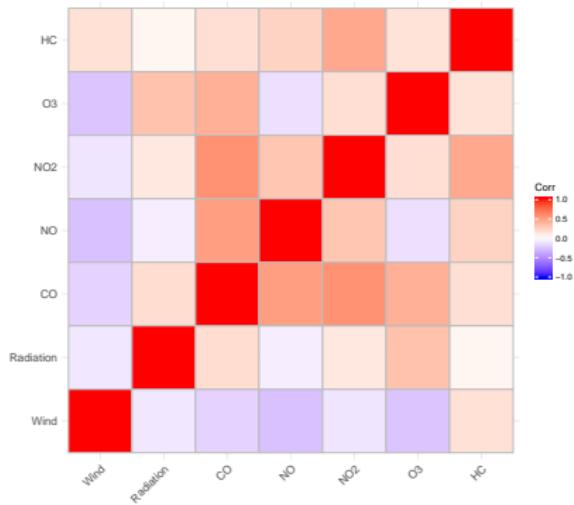
# Exploratory plots for multivariate data

See [Lecture1\\_2.R](#) for covariance matrix heatmap examples.



# Exploratory plots for multivariate data

See [Lecture1\\_2.R](#) for covariance matrix heatmap examples.



## Homework exercise 3

There are many different R functions for data visualization. Two R libraries that are particularly useful for multivariate data are the `lattice` library and the `ggplot2` library.

A good book on the use of lattice is Seepayan Sarkar, "Lattice. Multivariate Visualization with R", Springer, 2008. There are good blogs on the web that give you the `ggplot2` functions that are equivalent to the lattice functions, for example, <https://learnr.wordpress.com/2009/06/28/ggplot2-version-of-figures-in-lattice-multivariate-data-visualization-with-r-part-1/>.

Use R to visualize the data on National Track records for Men in a creative manner. Prepare a presentation of your visualization that shows both the graphs and the R-code that generated the graphs, as well as an interpretation of what the graphs are telling you about the variables and the associations in the data.

## Looking forward: Dimension reduction

So far we have concentrated on pairwise associations of the many variables or on at most 3-dimensional displays. But how do we visualize multivariate data with many more variables, or equivalently, much higher dimensions?

This is going to be a topic of much discussion later in this course.