

Index of Lecture 7: Introduction to multivariate analysis

Page	Title
1	Practical information
2	Textbook(s)
3	Multivariate statistical analysis
4	Example M1.1: sparrows
5	Example M1.2: skulls
6	Example M1.3: butterflies
7	Example M1.4: prehistoric dogs
8	Example M1.5: employment
9	Common assumptions
10	Matrix definitions
11	Matrix operations
12	Matrix inversion and quadratic form
13	Mean vectors and covariance matrices
14	Eigenvalues and eigenvectors
15	Graphing multivariate data
16	Sparrow: matrix plot
17	Dog: Chernoff plot
18	Dog: Profile plot
19	Structuring multivariate data

PRACTICAL INFORMATION

Today's lecture — preliminaries for multivariate part of course:

- introduction to **five datasets** used throughout, to illustrate possible objectives of multivariate analysis (Chapter 1 of Manly¹, in short M.1),
- gentle introduction to **matrix notation and algebra** (Chapter M.2) — very helpful to understand methods even in a non-technical sense (**our aim**),
- some **graphical tools** to display multivariate data (Chapter M.3).

Aim of course coverage of multivariate methods, adapted from Manly's book:

The purpose of this course is to introduce multivariate statistical methods to non-mathematicians. It is not intended to be a comprehensive course. Rather the intention is to keep the details to a minimum while serving as a practical guide that illustrates the possibilities of multivariate statistical methods.

Other news:

- **home assignment #2** due tomorrow (only students not in VHM 812),
- **need to discuss** format/logistics of classes during the Canada Games period (starting with tomorrow's lab),
- time to start thinking about/planning your **project**.

¹ The 4th edition of the book is by Manly & Alberto (2016), but Manly was sole author of the earlier editions.

TEXTBOOK(S)

Major textbook:

(M) Manly BFJ, Alberto JAN (2016), *Multivariate Statistical Methods: A Primer*, 4th ed.; only addition from 3rd ed. (2004) is R code (even a brief introduction to R).²

Supplementary texts (some suggestions, based on availability):

- (TF) Tabachnick BG, Fidell LS (2006/12/18), *Using Multivariate Statistics*, 5th/6th/7th ed.³; very thorough with well-worked examples, includes many topics not covered in course,
- (SL) James G, Witten D, Hastie T, Tibshirani R (2013/21), *An Introduction to Statistical Learning, with Applications in R*, 1st/2nd ed.⁴; excellent introduction to classification methods (in a broad sense),
- (ED) Everitt BS, Dunn G (2001), *Applied Multivariate Data Analysis*, 2nd ed.⁵; applied text focusing on classical methods.

Anticipated use of textbooks for course:

- (1) lectures will follow M closely, including using the provided datasets for illustration of methods,
- (2) other books or articles will be used, as needed and for methods not covered by these texts (e.g. PERMANOVA, SATSCAN) — feel free to contribute material of interest.

² The 1st edition from 1986 is at the Robertson Library; many parts are essentially unchanged (however, the 1st edition has no Chapter 3 on graphical methods).

³ The 5th edition is at the Robertson Library.

⁴ Freely downloadable from <https://www.statlearning.com/>.

⁵ Available at Robertson Library as E-book; also Everitt (2005), *An R and S-Plus Companion to Multivariate Analysis*.

MULTIVARIATE STATISTICAL ANALYSIS

The **terminology** is not definitive, but in our usage...

- does **not** include methods such as multiple regression with a single (“univariate”) outcome per subject but multiple predictor variables,⁶
- involves having multiple outcomes on each “subject”, viewed together as a multi-dimensional set of values (but there may also be (multiple) predictors).

Overview of (main) methods in course (others may be added):

- **multivariate regression and MANOVA**: ~ ordinary regression/ANOVA but for multivariate outcome,
- **principal components and factor analysis**: reduce the information in a set of variables to suitable new variables combined linearly,
- **classification and discriminant analysis**: separate existing groups in data based on suitable linear or non-linear combinations of the original variables,
- **cluster analysis**: identify groups in data based on suitable linear combinations of the original variables, or based on distances in a space for the observations,
- **multidimensional scaling**: create distances between observations along multiple directions.

⁶ Also termed **multivariable** models, to distinguish them from multivariate models.

EXAMPLE M1.1: SPARROWS

Data: 5 morphological measurements on 49 moribund female sparrows, of which 21 subsequently survived, taken to a biological laboratory after a storm (in 1898!),

- X_1 = total length (*mm*),
- X_2 = alar extent (*mm*),
- X_3 = length of beak and head (*mm*),
- X_4 = length of humerus (*mm*),
- X_5 = length of keel of sternum (*mm*),

Bird	X_1	X_2	X_3	X_4	X_5	survival
1	156	245	31.6	18.5	20.5	1
2	154	240	30.4	17.9	19.6	1
...		
48	162	245	32.5	18.5	21.1	0
49	164	248	32.3	18.8	20.9	0

Possible **objectives** of analysis:

- (a) describe relations between variables $X_1 - X_5$, within each of the survivor groups, and compare the two groups in terms of means, variability and relationships,
- (b) construct some combined feature(s) from $X_1 - X_5$ that allow(s) accurate prediction of survival, \sim index of fitness of the birds.

First thoughts:

- o descriptive analysis for (a): means, standard deviations, correlations,
 - o (b) could be approached by multiple logistic regression,
- \Rightarrow multivariate regression, principal components, classification/discriminant analysis.

EXAMPLE M1.2: SKULLS

Data: 4 morphometric measurements on 150 Egyptian male skulls, 30 from each of 5 distinct time periods,⁷

- X_1 = maximum breadth (*mm*),
- X_2 = basibregmatic height (*mm*),
- X_3 = basiolvelar length (*mm*),
- X_4 = nasal heightlength of humerus (*mm*),

Skull	X_1	X_2	X_3	X_4	period
1	131	138	89	49	1
2	125	131	92	48	1
...		
150	136	133	97	51	5

Possible **objectives** of analysis:

- (a) describe relations between $X_1 - X_4$ within each of the periods, and compare the 5 groups in terms of means, variability and relationships,
- (b) measure distances of distributions of $X_1 - X_4$ across the periods, and relate such distances to physical time,
- (c) construct a combined (linear) function from $X_1 - X_4$ that in some sense describes development over time.

First thoughts (beyond Example M1.1):

- not obvious how to quantify the impact of time (e.g., considering the 5 periods as ordered groups does not directly involve time differences)...

⇒ multivariate regression, multivariate distance, classification/discriminant analysis.

⁷ (1): early predynastic (4000 B.C.); (2): late predynastic (3300 B.C.); (3): 12-13th dynasties (1850 B.C.); (4): Ptolemaic (200 B.C.); (5): Roman (A.D. 150).

EXAMPLE M1.3: BUTTERFLIES

Data: 4 environmental and 6 genetic variables (frequencies for Phosphoglucose-Isomerase (Pgi) genes) for 16 colonies of the butterfly *Euphydryas editha* in California and Oregon (one colony only),

X_1 = altitude (*ft*),

X_2 = annual precipitation (*in*),

X_3 = maximum temperature ($^{\circ}\text{F}$),

X_4 = minimum temperature ($^{\circ}\text{F}$),

Y_r = frequency (%) for Pgi gene type r ; $r = 0.4, 0.6, 0.8, 1.0, 1.16, 1.3$.

Colony	X_1	X_2	X_3	X_4	$Y_{0.4}$	$Y_{0.6}$	$Y_{0.8}$	$Y_{1.0}$	$Y_{1.16}$	$Y_{1.3}$
SS	500	43	98	17	0	3	22	57	17	1
SB	808	20	92	32	0	16	20	38	13	13
...				
GL	10500	50	81	-12	0	3	1	92	4	0

Possible **objectives** of analysis:

- (a) describe relation between the Pgi gene frequencies and environmental variables,
- (b) is environmental similarity or physical proximity more important for genetic similarities? — the latter presumably representing effects of migration.

First thoughts:

- o dataset seems quite small, and contains no grouping of interest,
 - o the grouping of interest is among variables, not observations,
 - o spatial information needed for (b),
- \Rightarrow multivariate distances, canonical correlations.

EXAMPLE M1.4: PREHISTORIC DOGS

Data: 6 mandible (lower jaw) morphometric measurements (all in *mm*) on craniums of prehistoric dogs in Thailand and 6 other possibly related species,

X_1 = breadth of mandible,

X_2 = height of mandible
below 1st molar,

X_3 = length of 1st molar,

X_4 = breadth of 1st molar,

X_5 = length 1st molar – 3rd molar,

X_6 = length 1st molar – 4th premolar,

Species	X_1	X_2	X_3	X_4	X_5	X_6
Modern dog	9.7	21.0	19.4	7.7	32.0	36.5
Golden jackal	8.1	16.7	18.3	7.0	30.3	32.9
Chinese wolf	13.5	27.3	26.8	10.6	41.9	48.1
Indian wolf	11.5	24.3	24.5	9.3	40.0	44.6
Cuon	10.7	23.5	21.4	8.5	28.8	37.6
Dingo	9.6	22.6	21.1	8.3	34.4	43.1
Prehistoric dog	10.3	22.1	19.1	8.1	32.2	35.0

Possible **objectives** of analysis:

- (a) describe relations between the prehistoric dog and other species, in particular quantify distances between species.

First thoughts:

- dataset seems quite small (e.g., lacks replication within species), and contains no grouping of interest.

⇒ multivariate distances, cluster analysis.

EXAMPLE M1.5: EMPLOYMENT

Data: Percentages of workforce employment in 9 different sectors⁸ of 30 European countries grouped into 4 political/economical zones⁹, as of the 1990s (approximately),

X_r = employment (%)
in sector r ;
 $r = 1, \dots, 9$,

Country	Political group	X_1 AGR	X_2 MIN	X_3 MAN	X_4 PS	X_5 CON	X_6 SER	X_7 FIN	X_8 SPS	X_9 TC
Belgium	EU	2.6	0.2	20.8	0.8	6.3	16.9	8.7	36.9	6.8
Denmark	EU	5.6	0.1	20.4	0.7	6.4	14.5	9.1	36.3	7.0
...				
Turkey	Other	44.8	0.9	15.3	0.2	5.2	12.4	2.4	14.5	4.4

Possible **objectives** of analysis:

(a) describe employment patterns, and if possible establish groups of countries with similar patterns; such groupings could then be compared with the political grouping provided.

First thoughts:

- primary interest probably not related to the grouping provided,
 - it may pose problems that the employment percentages add up to 100%,
- ⇒ multivariate distances, principal components/factor analysis, classification/
discriminant analysis, cluster analysis.

⁸ (1): AGR (agriculture); (2): MIN (mining); (3): MAN (manufacturing); (4): PS (power and water supplies), (5): CON (construction); (6): SER (services); (7): FIN (finances); (8): SPS (social and personal services); (9): TC (transport and communications).

⁹ (1): EU (European union); 2: EFTA (European free trade area); 3: East (former Eastern Europe union); (4): Other.

COMMON ASSUMPTIONS

Two common assumptions:

- **independence** between sets of measures taken on different “subjects” (whereas independence is never assumed for the multiple measures on the same subject),
 - * can be violated by data possessing a particular structure, e.g. hierarchical,
 - * generalizations to dependent data depend on the specific methods but are generally not straightforward (\Rightarrow tempting to ignore and “interpret with caution”),
- **normally distributed, interval-scale measures**, more precisely that the p -dimensional set of measures, say X_1, \dots, X_p , follows a MVN (multivariate normal) distribution with mean μ and variance-covariance matrix Σ , where
 - * μ is the set of means for the components,
 - * Σ contains the covariances between all pairs of components (see also 7L–13);this assumption includes a **normal distribution for each component**¹⁰; violations of normality can be dealt with in different ways,
 - * transformation to “better” scale (e.g. for skewed distributions),
 - * specialized methods, such as correspondence analysis for categorical data,
 - * claim of robustness of methods towards non-normality (to be further discussed).

¹⁰ A MVN also requires the conditional distributions of subsets of components given the others (or all linear combinations of the components) to be normal, but this extra condition is often ignored in practice; M1.3, TF (p. 78).

MATRIX DEFINITIONS

An $m \times n$ matrix A is a collection of mn values organized in m rows and n columns:

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ \vdots & & & \ddots & \\ a_{m1} & a_{m2} & a_{m3} & \dots & a_{mn} \end{pmatrix}$$

- a **square matrix** has equally many rows and columns: $m = n$,
- a **row vector**, say $r = (r_1 \ r_2 \ \dots \ r_n)$, has only one row: $m = 1$,
- a **column vector** has only one column: $n = 1$,
- the **transpose** matrix A' or A^t results from interchanging the rows and columns (i.e., $a_{ij}^t = a_{ji}$),
 - * a **symmetric** matrix A satisfies: $A^t = A$,
 - * the transpose of a row vector is a column vector, and vice versa,
- a **zero** matrix 0 has all elements $= 0$; a **one** matrix 1 has all elements $= 1$,
- a square **diagonal** matrix D has only 0s outside the diagonal,¹¹
- the **unity** matrix I is a diagonal matrix with 1s in the diagonal,¹²
- the **trace** of a square matrix is the sum of the diagonal elements,

$$\text{tr}(A) = a_{11} + a_{22} + \dots + a_{nn}.$$

¹¹ The condition can be expressed that $d_{ij} = 0$ when $i \neq j$.

¹² The condition can be expressed that $i_{ij} = 1$ when $i = j$, and $i_{ij} = 0$ otherwise.

MATRIX OPERATIONS

Let A and B be matrices, and let k be a (real-valued) number, often termed a **scalar**, so as to distinguish it from vectors and matrices.

- **addition** ($A + B$) and **subtraction** ($A - B$) of matrices occurs for each element (at its position in the matrix),
 - * A and B must have the same dimensions,
 - * e.g.: $(A + B)_{11} = a_{11} + b_{11}$; $(A - B)_{32} = a_{32} - b_{32}$,
- **scalar multiplication** by k also occurs for each element, e.g. $(kA)_{11} = k \cdot a_{11}$,
- **matrix multiplication** ($A \cdot B$ or $A.B$ or $A \times B$, but most commonly just AB) is **not** defined by multiplying pairs of matrix elements,
 - * requires **compatible matrix dimensions**: no. columns of A must equal no. of rows of B , say A ($m \times n$) and B ($n \times p$),
 - * for A ($m \times n$) and B ($n \times p$), $A.B$ has dimension ($m \times p$),
 - * the (i, k) th element of $A.B$ is the sum

$$(A.B)_{ik} = \sum_j a_{ij} b_{jk} = a_{i1}b_{1k} + a_{i2}b_{2k} + \dots + a_{in}b_{nk},$$

- * it is **not** generally true that $A.B = B.A$; note that $A.B$ and $B.A$ are not even of the same dimensions, unless A, B are square matrices.

MATRIX INVERSION AND QUADRATIC FORM

Matrix inversion = the matrix version of **division** but more complicated than for numbers.

Definition: a square matrix A ($n \times n$) has inverse A^{-1} ($n \times n$) if the matrices satisfy:

$$A \cdot A^{-1} = A^{-1} \cdot A = I.$$

Example: the 2×2 matrix inverse of $A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$

◦ exists exactly when the **determinant**¹³ $\Delta = a_{11}a_{22} - a_{21}a_{12}$ is nonzero,

◦ can be computed as $A^{-1} = \begin{pmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{pmatrix} / \Delta$.

General definitions and results for A ($n \times n$):

◦ A has an inverse A^{-1} exactly when the determinant¹³ $|A|$ is non-zero,

◦ if A^{-1} does not exist, then A is called **singular**,

◦ if A^{-1} exists and $A^{-1} = A^t$, then A is called **orthogonal**.

A **quadratic form** Q for a symmetric matrix A ($n \times n$) is a linear function of a column vector $\mathbf{x} = (x_1 \ x_2 \ \dots \ x_n)^t$,

$$Q = Q(\mathbf{x}) = \mathbf{x}^t A \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j.$$

¹³ The standard notation for the determinant of A is $|A|$. Also for square matrices larger than (2×2), the determinant $|A|$ has an explicit, but complex form; $|A|$ is therefore best calculated using mathematical/statistical software.

MEAN VECTORS AND COVARIANCE MATRICES

Multivariate normal (MVN) distribution for a p -dimensional column vector $\mathbf{X} = (X_1 \ X_2 \ \dots \ X_p)^t \sim \text{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where

- * $\boldsymbol{\mu} = (\mu_1 \ \mu_2 \ \dots \ \mu_p)^t$ is a p -dimensional column vector of means (i.e., $\text{E}X_j = \mu_j$),
- * $\boldsymbol{\Sigma} = (\sigma_{jk})_{jk}$ is a $(p \times p)$ **positive definite**¹⁴ matrix of the pairwise covariances,¹⁵

has p -dimensional probability density function

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

Estimation of mean and covariance:

From i.i.d.¹⁶ multivariate observations X_1, X_2, \dots, X_n we can estimate the mean vector and covariance by

- * $\hat{\boldsymbol{\mu}} = \bar{\mathbf{X}} = (\bar{X}_1 \ \bar{X}_2 \ \dots \ \bar{X}_p)^t$, the column vector of means for each variable,
- * $\hat{\boldsymbol{\Sigma}} = \mathbf{S} = (s_{jk})_{jk}$, the **empirical covariance matrix** with the $s_{jk} = r_{jk} s_j s_k$ defined in terms of the Pearson correlation coefficient (r_{jk}) between X_j and X_k and the standard deviations (s_j, s_k) for X_j and X_k , respectively.

¹⁴ A symmetric matrix \mathbf{A} is positive definite if: $Q(\mathbf{x}) = \mathbf{x}^t \mathbf{A} \mathbf{x} > 0$ for any non-zero \mathbf{x} ; see also next slide.

¹⁵ Formally, $\sigma_{jk} = \text{E}(X_j X_k) - (\text{E}X_j)(\text{E}X_k)$, and $\sigma_{jj} = \text{Var}(X_j)$; also, the correlation matrix $\mathbf{R} = (\rho_{jk})_{jk}$ gives the off-diagonal correlations $\rho_{jk} = \sigma_{jk} / (\sigma_{jj} \sigma_{kk})$.

¹⁶ Also here, i.i.d. stands for: independent and identically distributed observations.

EIGENVALUES AND EIGENVECTORS

For a symmetric matrix \mathbf{A} ($n \times n$), a scalar λ and a column vector \mathbf{x} ($n \times 1$), we consider the equation

$$\mathbf{Ax} = \lambda \mathbf{x} \quad \text{or} \quad (\mathbf{A} - \lambda \mathbf{I}) \mathbf{x} = \mathbf{0},$$

which can also be represented as the set of equations:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = \lambda x_1, \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = \lambda x_2, \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = \lambda x_n. \end{cases}$$

The possible solutions¹⁷ reflect fundamental properties of \mathbf{A} :

- the **eigenvalues** (or characteristic roots) λ , of which there can be at most n distinct values, all of which are **> 0 for a positive definite matrix**, and which satisfy

$$\text{tr}(\mathbf{A}) = \lambda_1 + \dots + \lambda_n,$$

- * the **eigenvectors** \mathbf{x} for each λ express specific “directions” in the matrix.

Applying these general matrix concepts to a covariance matrix (Σ) gives:

- $\text{tr}(\Sigma) = \sigma_{11} + \dots + \sigma_{pp} = \lambda_1 + \dots + \lambda_n$,
— i.e., the eigenvalues decompose the total variance...
- the eigenvectors correspond to specific directions of the variation in the data.

¹⁷ Calculation of solutions λ and \mathbf{x} to the eigenvalue problem is not straightforward to establish; TF work through a (2×2) example.

GRAPHING MULTIVARIATE DATA

Fundamental challenge: getting beyond 2-d (or perhaps 3-d) scatterplots to display relationships between variables.

Extension to 3-d plots:

- use 3rd axis (depths), and rotate points,
- or create (interpolated/smooth) surfaces for 3rd variable,
- + both ideas useful, but not exactly grand extensions. . .

First idea: pairwise scatter plots \Rightarrow **matrix** or draftman's plot,

- + can show relations between a limited number of variables,
- difficult to see similarity between observations.

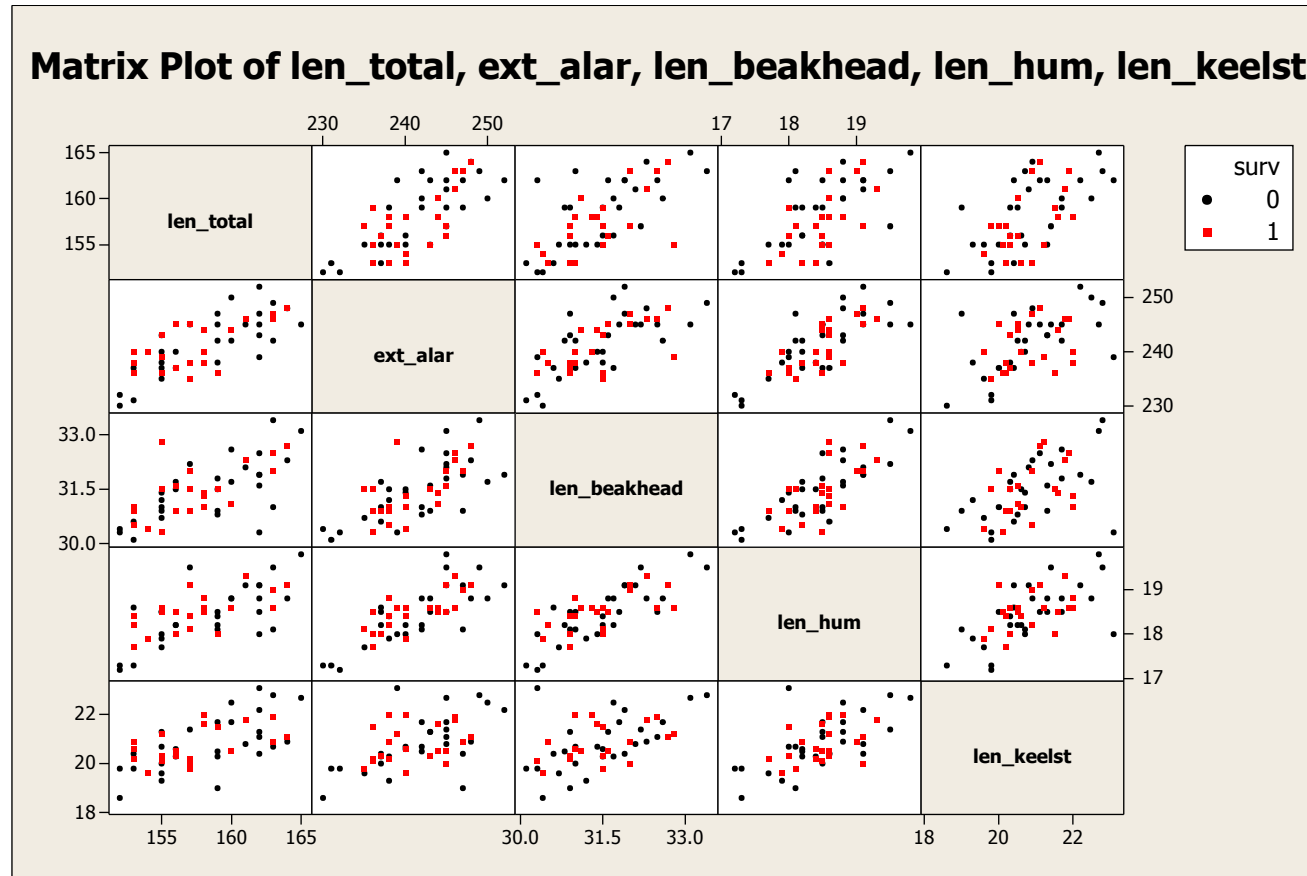
Next idea: represent variables' values symbolically in displays involving multiple features, e.g. faces (**Chernoff**) or stars.

Another idea: position multiple measures per subject on x -axis at suitable (arbitrary) points and join values for same subject by line (**profile** plot), or position them as multiple bars beside each other (**bar** graph),

- + can illustrate grouping and distance between subjects,
- requires measures on comparable scales, difficult to see relation between variables.

Final comment: one key purpose of multivariate methods is to **reduce the dimension of information** across many variables into a few (“index”; Manly) variables, e.g. for the purpose of facilitating graphical exploration.

SPARROW: MATRIX PLOT

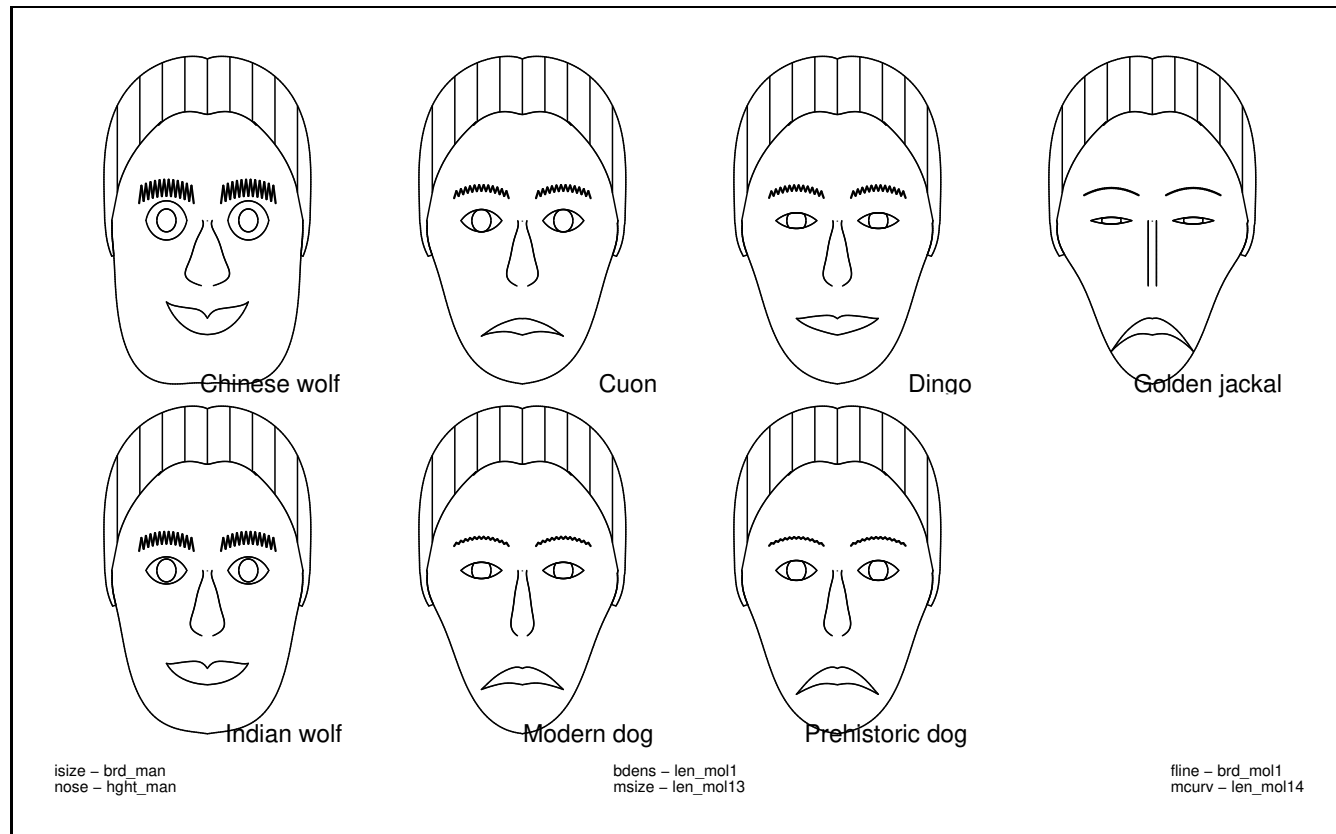


- all variables seem positively and approximately linearly associated, to similar degree (correlation range: 0.526 – 0.769),
- no obvious differences between survivors and mortals,
- no apparent strong outliers.

DOG: CHERNOFF PLOT

Idea¹⁸: use different facial features (up to 18) to represent variables,

- one selected feature per variable, scaled so that low/high values correspond to extreme features,
- some arbitrariness/flexibility in the allocation of features to variables,
- not clear that visual impression is not affected substantially by allocation of features to variables (other symbols, e.g. stars, may be more robust).



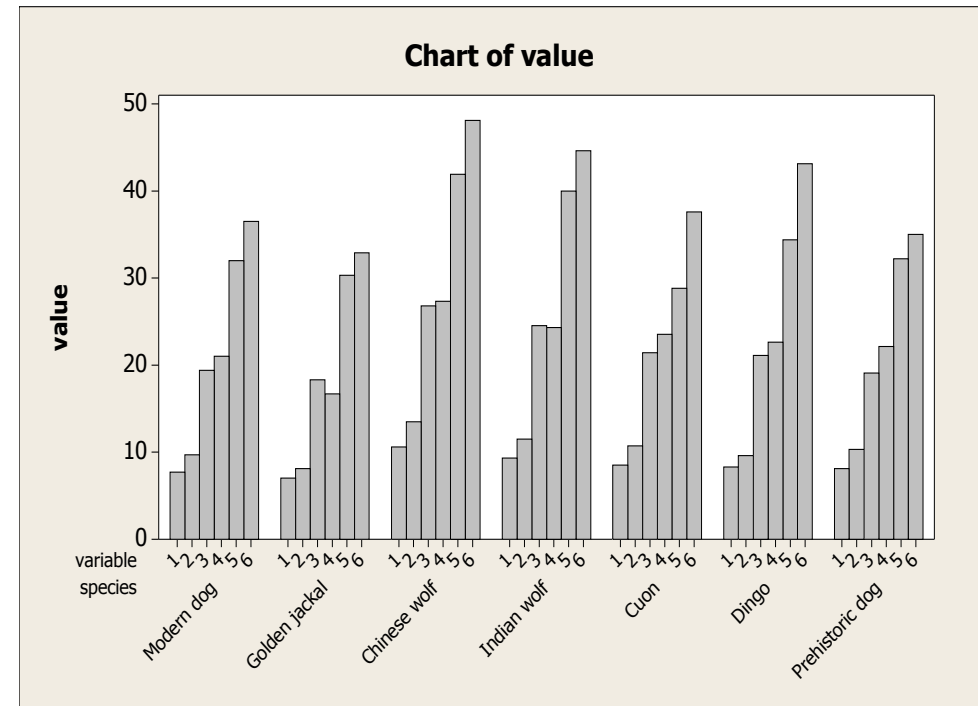
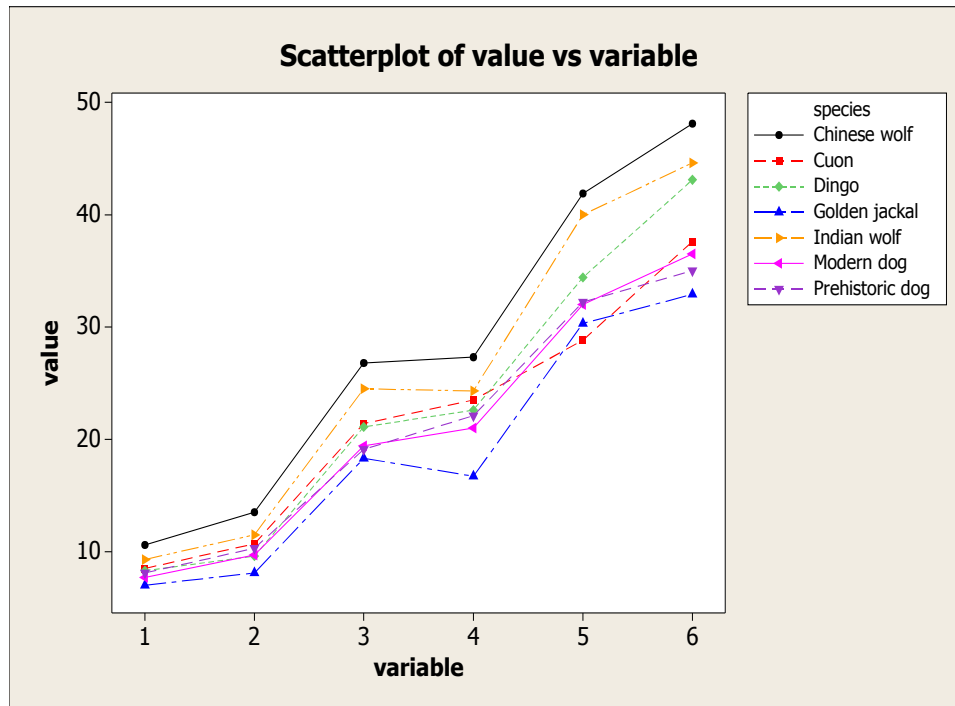
Interpretation:

- the plot¹⁹ shows closeness of prehistoric and modern dog, also their distance to the wolves.

¹⁸ Attributed to Chernoff H (1973), *J. Amer. Statist. Assoc.* 68, 361-368.

¹⁹ Done with Stata chernoff add-on package, using the features: isize (eye size), nose (nose line), bdens (brow density), msize (mouth size), mcurv (mouth curvature) and fline (face line).

DOG: PROFILE PLOT AND BAR GRAPH



Interpretation:

- both plots show closeness of prehistoric and modern dog, and their differences from the wolves (Chinese and Indian).

STRUCTURING MULTIVARIATE DATA

Two major **formats of datasets**:

- **wide**: usual format, one record (worksheet row) per subject, with multiple columns for the multivariate measurements as well as any extra (predictor) variables,²⁰
- **long**: one record per measurement, with the multivariate measurements put together in a single column and the type of measurement indicated in another column, with one additional column containing subject id's, as well as (repeated values of) predictor variables in other columns,
 - * most common data format for univariate data and also for multivariate data, when the multivariate measurements are different instances of the same variable, e.g. over time (i.e., repeated measures data).

How to **switch** between data formats? — software specific implementations:²¹

- **Minitab**: “stacking” columns for: wide → long; “unstacking” for reverse process,
- **Stata**: “reshape” command in versions “wide” and “long”, indicating the targeted data format,

Why shift data format? — typically due to specific requirements of implementations of analytical or graphical procedures.

²⁰ All examples of the lecture showed the data in this format.

²¹ Note also the related operation **transpose** (Stata: “xpose”), switching rows and columns in a worksheet.