

UNIVERSITA' DEGLI STUDI DI PADOVA

FACOLTA' DI SCIENZE STATISTICHE

CORSO DI LAUREA IN STATISTICA, ECONOMIA E FINANZA



TESI DI LAUREA

**LA REGRESSIONE MULTIVARIATA:
UN'APPLICAZIONE A DATI ANTROPOLOGICI**

Relatore: Ch.mo Prof. Laura Ventura

Laureanda: Marina Marzovilli

Matricola: 515937

ANNO ACCADEMICO 2006-2007

AI MIEI GENITORI

Indice

Introduzione	1
1 Inferenza nel modello di regressione multivariata	3
1.1 Stima di massima verosimiglianza dei coefficienti di regressione e della matrice Σ	3
1.2 La verifica d'ipotesi	6
1.2.1 LRT	6
1.2.2 UIT	7
1.3 Intervalli di confidenza simultanei per \mathbf{B}	9
1.4 Alcune considerazioni conclusive	9
1.4.1 Stima dei coefficienti nel caso in cui \mathbf{X} non sia di rango pieno	9
1.4.2 Valutazione dell'effetto della media sulle risposte	10
1.4.3 Dalle stime OLS alle stime GLS	11
2 MANOVA: Analisi della varianza multivariata	13
2.1 Test multivariati	14
2.1.1 T^2 di Hotelling	14
2.1.2 Lambda di Wilks	15
2.1.3 Pillai-Bartlett <i>trace</i>	16
2.1.4 GCR	16

2.1.5 Alcune conclusioni	16
2.2 Test “ <i>post-hoc</i> ”	18
2.3 Assunzioni alla base dell’analisi multipla della varianza	19
2.4 Utilizzi della MANOVA	21
2.5 Conclusioni	21
3 Un caso studio	23
3.1 I dati	23
3.2 Analisi grafiche univariate	27
3.3 Analisi bivariate	33
3.4 Conclusioni	38
4 Adattamento ai dati	43
4.1 Il pacchetto <i>lm</i>	43
4.1.1 Il comando <code>anova.mlm()</code>	44
4.2 La stima del modello completo	45
4.2.1 Semplificazione del modello	47
4.3 Previsioni	50
4.4 Conclusioni	53
Appendice	57
Bibliografia	61

Introduzione

La regressione multivariata riguarda il problema della modellazione di più di una variabile risposta a partire dallo stesso insieme di variabili esplicative. È possibile, in questo caso, cercare di capire come si può tener conto dalle correlazioni esistenti tra le risposte per comparare l'accuratezza della previsione con quella delle procedure che generalmente si usano nel caso di regressioni univariate, nelle quali per ogni singola risposta viene utilizzato un insieme di regressori.

Nel presentare il caso in cui la risposta consiste in $p > 1$ variabili differenti si possono facilmente estendere le considerazioni fatte per la regressione multipla all'analisi multivariata. La regressione multivariata è definita dal modello

$$\mathbf{Y} = \mathbf{XB} + \mathbf{U}, \quad (1)$$

dove \mathbf{Y} non è più il vettore ($n \times 1$) considerato nel caso dell'analisi univariata, con n che rappresenta il numero di dati osservati, bensì una matrice ($n \times p$) contenente p variabili risposta; \mathbf{X} è una matrice ($n \times q$) di variabili esplicative (con $q \geq 1$); \mathbf{B} è una matrice ($q \times p$) di coefficienti di regressione e \mathbf{U} è la matrice dei termini d'errore di dimensione ($n \times p$).

Il modello (1) può essere scritto come

$$\begin{bmatrix} y_{11} & y_{12} & \cdots & y_{1p} \\ y_{21} & y_{22} & \cdots & y_{2p} \\ \cdots & \cdots & \cdots & \cdots \\ y_{n1} & y_{n2} & \cdots & y_{np} \end{bmatrix} = \begin{bmatrix} 1 & x_{12} & \cdots & x_{1q} \\ 1 & x_{22} & \cdots & x_{2q} \\ \cdots & \cdots & \cdots & \cdots \\ 1 & x_{n2} & \cdots & x_{nq} \end{bmatrix} \begin{bmatrix} \beta_{01} & \beta_{02} & \cdots & \beta_{0p} \\ \beta_{11} & \beta_{12} & \cdots & \beta_{1p} \\ \cdots & \cdots & \cdots & \cdots \\ \beta_{q1} & \beta_{q2} & \cdots & \beta_{qp} \end{bmatrix} + \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1p} \\ u_{21} & u_{22} & \cdots & u_{2p} \\ \cdots & \cdots & \cdots & \cdots \\ u_{n1} & u_{n2} & \cdots & u_{np} \end{bmatrix}$$

Siano \mathbf{y}_i e \mathbf{x}_i l' i -esima ($i = 1, \dots, n$) riga, rispettivamente, delle matrici \mathbf{Y} e \mathbf{X} ; nel caso del modello con intercetta, allora la prima colonna della matrice \mathbf{X} è il vettore unitario. Se \mathbf{X} è una matrice di q variabili esplicative osservate su ognuno degli n individui del campione considerato, allora parliamo di *modello di regressione multivariata*. Nel modello (1) le colonne di \mathbf{Y} , che rappresentano le variabili dipendenti del modello osservate sugli n individui del campione, sono spiegate dalle colonne di \mathbf{X} . In particolare vale la relazione

$$E(\mathbf{y}_{ij}) = \mathbf{x}_i^T \mathbf{B}_j, \quad i = 1, \dots, n, j = 1, \dots, p,$$

ossia il valore atteso di \mathbf{y}_{ij} dipende dall' i -esima riga della matrice \mathbf{X} e dalla j -esima colonna \mathbf{B}_j della matrice dei coefficienti di regressione \mathbf{B} . Nel caso in cui $p = 1$, abbiamo di fronte il modello di regressione multipla, con una sola variabile risposta, per il quale possiamo scrivere nel modo usuale $\mathbf{y} = \mathbf{X}\mathbf{B} + \boldsymbol{\varepsilon}$.

Nella maggior parte dei casi è possibile assumere che la matrice \mathbf{U} degli errori sia distribuita normalmente: \mathbf{U} è una matrice con una distribuzione $N_p(0, \boldsymbol{\Sigma})$, dove $\boldsymbol{\Sigma}$ è la matrice di covarianza. Si assume, inoltre, che la distribuzione di \mathbf{U} sia incorrelata con \mathbf{X} , ossia che $E(\mathbf{U}|\mathbf{X}) = \mathbf{0}$. Sotto l'assunzione di normalità degli errori, la log-verosimiglianza per i parametri \mathbf{B} e $\boldsymbol{\Sigma}$ è data da

$$l(\mathbf{B}, \boldsymbol{\Sigma}) = -\frac{1}{2} n \log |2\pi \boldsymbol{\Sigma}| - \frac{1}{2} \text{tr}[(\mathbf{Y} - \mathbf{X}\mathbf{B}) \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X}\mathbf{B})^T], \quad (2)$$

dove il simbolo $\text{tr}(\mathbf{A})$ indica la traccia della matrice \mathbf{A} , ossia la somma degli elementi che si trovano sulla diagonale di tale matrice.

Nei capitoli successivi verranno sviluppati i seguenti argomenti: il primo capitolo presenterà alcuni concetti generali relativi alla teoria della regressione multivariata e all'inferenza sui parametri stimati nel modello oggetto di studio; nel secondo capitolo si tratterà dell'analisi multivariata della varianza (MANOVA) e dei test statistici ad essa correlati; nel terzo capitolo saranno presentati i dati oggetto di studio e verranno eseguite analisi grafiche preliminari e test su di essi; infine, nell'ultimo capitolo si stimerà il modello di regressione multivariata sui dati, valutando gli effetti che hanno le variabili esplicative sulle risposte considerate; in appendice sono contenute alcune informazioni relative alle distribuzioni multivariate legate ai test statistici presentati nel seguito.

Capitolo 1

Inferenza nel modello di regressione multivariata

Nel modello (1), per la stima dei parametri \mathbf{B} e Σ è possibile utilizzare il metodo della massima verosimiglianza, che si basa sulla soluzione di equazioni di stima ottenute a partire dalla (2). Si ottengono così gli stimatori di massima verosimiglianza per la matrice dei coefficienti \mathbf{B} e per la matrice Σ di varianza e covarianza degli errori, dove Σ è una matrice definita positiva e simmetrica, di dimensioni $(p \times p)$. A partire dai risultati ottenuti è possibile eseguire verifiche di ipotesi sui parametri del modello e costruire intervalli di confidenza. È possibile, inoltre, considerare delle generalizzazioni del modello quando alcune delle ipotesi considerate non valgono. Ad esempio, si può vedere come cambiano le stime dei coefficienti nel caso in cui i termini d'errore siano eteroschedastici, ossia con varianze diverse.

1.1 Stima di massima verosimiglianza dei coefficienti di regressione e della matrice Σ

Sia $\mathbf{P} = \mathbf{I} - \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$, dove \mathbf{I} è la matrice identità di ordine n . La matrice \mathbf{P} gode di alcune proprietà: è simmetrica ed idempotente ($\mathbf{P} = \mathbf{P}^T$, $\mathbf{P} \times \mathbf{P} = \mathbf{P}^2 = \mathbf{P}$) ed ha rango $n-q$. Essa rappresenta la proiezione ortogonale in \mathbf{R}^n delle colonne di \mathbf{X} . In particolare, $\mathbf{P}\mathbf{X} = \mathbf{0}$.

Nel modello (1) considerato si ipotizza che le righe della matrice \mathbf{U} siano indipendenti e distribuite normalmente con un vettore p -dimensionale di medie pari a zero e con una

matrice di varianza e covarianza Σ ($p \times p$). Ciò implica che per ogni colonna della matrice delle risposte ci sia una regressione completa col suo insieme di coefficienti di regressione. La matrice stimata dei coefficienti è ottenuta come

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y},$$

dove la matrice $\hat{\mathbf{B}}$ è distribuita come una t di Student multivariata (si veda l'Appendice) con media \mathbf{B} e matrice di varianza $[(n - q - 1)(\mathbf{X}^T \mathbf{X})]^{-1} \mathbf{G}$, ossia

$$\hat{\mathbf{B}} \sim t(n - q - 1, \mathbf{B}, [(n - q - 1)(\mathbf{X}^T \mathbf{X})]^{-1} \mathbf{G})$$

con $\mathbf{G} = (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})^T (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})$.

In particolare $\hat{\mathbf{B}}_k$, che rappresenta la k -esima colonna di $\hat{\mathbf{B}}$ ($k = 1, \dots, p$), si distribuisce come una t di Student multivariata con media \mathbf{B}_k e matrice di varianza e covarianza $(n - q - p)^{-1} \mathbf{W}_{kk} \mathbf{G}$; ossia

$$\hat{\mathbf{B}}_k \sim t(n - q - p, \mathbf{B}_k, (n - q - p)^{-1} \mathbf{W}_{kk} \mathbf{G}),$$

dove \mathbf{W}_{kk} è l'elemento di posto (k, k) della matrice $\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1}$.

Inoltre, per il singolo stimatore $\hat{\mathbf{B}}_{jk}$ ($k = 1, \dots, p, j = 1, \dots, q$), la distribuzione di riferimento è

$$\hat{\mathbf{B}}_{jk} \sim t(n - q - p, \mathbf{B}_{jk}, (n - q - p)^{-1} \mathbf{W}_{kk} \mathbf{g}_{jj}),$$

che rappresenta una t di Student univariata con media \mathbf{B}_{jk} e varianza $(n - q - p)^{-1} \mathbf{W}_{kk} \mathbf{g}_{jj}$, dove \mathbf{g}_{jj} è l'elemento j -esimo sulla diagonale di \mathbf{G} . (si veda Rowe [6])

La stima della matrice dei coefficienti di regressione $\hat{\mathbf{B}}$ non dipende dalla covarianza rappresentata dalla matrice Σ o dalla sua stima; rimane la stessa a prescindere dalla presenza o meno di correlazione.

Per il modello (1), $\hat{\mathbf{B}}$ è una stima non distorta di \mathbf{B} . Infatti sostituendo $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{U}$ nell'equazione di $\hat{\mathbf{B}}$ si ottiene

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\mathbf{B} + \mathbf{U}) = \mathbf{B} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{U}$$

e quindi vale la relazione $E(\hat{\mathbf{B}}) = \mathbf{B}$ poichè $E(\mathbf{U}) = \mathbf{0}$. Infatti \mathbf{U} è la matrice dei termini di errore tale per cui le medie sono nulle e quindi $E(\hat{\mathbf{B}}) = \mathbf{B}$.

La covarianza tra $\hat{\mathbf{B}}_{jk}$ e $\hat{\mathbf{B}}_{kl}$ è data da $\sigma_{jl}w_{ik}$, dove σ_{jl} è l'elemento di posto (j,l) della matrice di varianza e covarianza degli errori Σ , mentre w_{ik} è l'elemento di posto (i,k) della matrice $\mathbf{W} = (\mathbf{X}^T\mathbf{X})^{-1}$. Infatti dall'espressione di $\hat{\mathbf{B}}$, si ricava che $\hat{\beta}_{ij} - \beta_{ij} = \mathbf{a}_i^T \mathbf{u}_j$, dove \mathbf{a}_i è la i -esima riga della matrice $\mathbf{A} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$, con $\mathbf{A}^T\mathbf{A} = \mathbf{W}$, e \mathbf{u}_j è la j -esima colonna della matrice degli errori \mathbf{U} . Allora

$$\text{Cov}(\hat{\beta}_{ij}, \hat{\beta}_{kl}) = \mathbf{a}_i^T \mathbf{E}[\mathbf{u}_i \mathbf{u}_l^T] \mathbf{a}_k = \mathbf{a}_i^T [\sigma_{jl} \mathbf{I}] \mathbf{a}_k = \sigma_{jl} w_{ik}. \quad (1.1)$$

Oltre ai coefficienti di regressione, per il modello si ottiene anche una stima della matrice Σ , come

$$\hat{\Sigma} = n^{-1} \mathbf{Y}^T \mathbf{P} \mathbf{Y}.$$

Si può dimostrare che per $\hat{\Sigma}$ vale

$$n \hat{\Sigma} \sim W_p(\Sigma, n - q),$$

con W_p che indica una distribuzione di Wishart con $n - q$ gradi di libertà e matrice di scala Σ . Tale distribuzione rappresenta una generalizzazione della distribuzione χ^2 nel caso multivariato (si veda l'Appendice).

Una volta calcolati i coefficienti di regressione, è possibile ottenere i valori predetti dal modello ("fitted values") come

$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\mathbf{B}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

Inoltre, i residui si calcolano nel seguente modo

$$\hat{\mathbf{U}} = \mathbf{Y} - \mathbf{X} \hat{\mathbf{B}} = \mathbf{P} \mathbf{Y}.$$

Da ciò si ottiene che $\mathbf{E}(\hat{\mathbf{U}}) = \mathbf{0}$. Infatti

$$\mathbf{E}(\hat{\mathbf{U}}) = \mathbf{E}(\mathbf{Y}) - \mathbf{E}(\mathbf{X} \hat{\mathbf{B}}) = \mathbf{X} \mathbf{B} - \mathbf{X} \mathbf{B} = \mathbf{0}.$$

Nel modello di regressione multivariata può essere necessario, come nel caso univariato, effettuare una selezione delle variabili esplicative da utilizzare come regressori. Nella selezione dei regressori bisogna tener conto del fatto che la matrice $\mathbf{X}^T \mathbf{X}$ deve essere invertibile e che è necessario scegliere solo le variabili risposta non troppo strettamente correlate tra loro. In particolare, per la selezione dei regressori, le procedure più utilizzate sono la selezione *backward* e quella *forward* che permettono in modo quasi automatico di eliminare o aggiungere le esplicative che non influenzano, o

influenzano, significativamente le variabili risposta. Per approfondimenti si veda Mardia *et al.* (1979).

1.2 La verifica d'ipotesi

Nel modello (1) si possono considerare ipotesi nulla della forma $H_0: \mathbf{C}_I \mathbf{B} \mathbf{M}_I = \mathbf{D}$, dove \mathbf{C}_I è una matrice ($g \times q$) di rango g , \mathbf{M}_I è una matrice ($p \times r$) di rango r e \mathbf{D} è una matrice ($g \times r$). In molti casi si può supporre che $\mathbf{D} = \mathbf{0}$ e $\mathbf{M}_I = \mathbf{I}$, e l'ipotesi nulla diventa semplicemente $H_0: \mathbf{C}_I \mathbf{B} = \mathbf{0}$. Le righe di \mathbf{C}_I permettono di capire qual è l'effetto sulla regressione di una combinazione lineare delle variabili indipendenti, mentre le colonne di \mathbf{M}_I focalizzano l'attenzione su particolari combinazioni lineari di variabili dipendenti. Ci sono due metodi per verificare tali ipotesi: il test del rapporto di verosimiglianza (*likelihood ratio test*, LRT) ed il test unione ed intersezione (*union intersection test*, UIT).

1.2.1 LRT

Considerando l'ipotesi nulla $H_0: \mathbf{C}_I \mathbf{B} = \mathbf{D}$ (con $\mathbf{M}_I = \mathbf{I}$), è opportuno definire ulteriori matrici per costruire il test. Sia \mathbf{C}_2 la matrice tale per cui $\mathbf{C}^T = (\mathbf{C}_1^T, \mathbf{C}_2^T)$ sia di ordine ($q \times q$) non singolare e sia \mathbf{B}_0 la matrice di dimensioni ($q \times p$) che soddisfa la relazione $\mathbf{C}_I \mathbf{B}_0 = \mathbf{D}$. È allora possibile riscrivere il modello (1) come

$$\mathbf{Y}_+ = \mathbf{Z} \mathbf{A} + \mathbf{U},$$

dove $\mathbf{Y}_+ = \mathbf{Y} - \mathbf{X} \mathbf{B}_0$, $\mathbf{Z} = \mathbf{X} \mathbf{C}^{-1}$ e $\mathbf{A} = (\mathbf{A}_1^T, \mathbf{A}_2^T) = \mathbf{C}(\mathbf{B} - \mathbf{B}_0)$. Allora l'ipotesi nulla $H_0: \mathbf{C}_I \mathbf{B} = \mathbf{D}$ può essere riscritta come $H_0: \mathbf{A}_1 = \mathbf{0}$. Considerando la partizione $\mathbf{C}^{-1} = (\mathbf{C}^{(1)}, \mathbf{C}^{(2)})$ si può definire la matrice di proiezione \mathbf{P}_I sul sottospazio ortogonale della colonne di $\mathbf{X} \mathbf{C}^{(2)}$, come

$$\mathbf{P}_I = \mathbf{I} - \mathbf{X} \mathbf{C}^{(2)} (\mathbf{C}^{(2)} \mathbf{X}^T \mathbf{X} \mathbf{C}^{(2)})^{-1} \mathbf{C}^{(2)T} \mathbf{X}^T.$$

Dall'espressione della funzione di verosimiglianza (2), si ottiene che sotto le ipotesi nulla e alternativa, le verosimiglianze massimizzate corrispondono rispettivamente a

$$|2\pi n^{-1} \mathbf{Y}_+^T \mathbf{P}_I \mathbf{Y}_+|^{-n/2} \exp\left\{-\frac{1}{2} n \mathbf{p}\right\} \quad \text{e} \quad |2\pi n^{-1} \mathbf{Y}^T \mathbf{P} \mathbf{Y}|^{-n/2} \exp\left\{-\frac{1}{2} n \mathbf{p}\right\}.$$

La statistica LRT si ottiene quindi come

$$\text{LRT} = |\mathbf{Y}^T \mathbf{P} \mathbf{Y}| / |\mathbf{Y}_+ \mathbf{P}_1 \mathbf{Y}_+|.$$

Se si definisce $\mathbf{P}_2 = \mathbf{P}_1 - \mathbf{P}$, anche \mathbf{P}_2 è una matrice di proiezione. In particolare, l'espressione di tale matrice sarà

$$\mathbf{P}_2 = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}_1^T [\mathbf{C}_1 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}_1^T]^{-1} \mathbf{C}_1 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T.$$

Inoltre, poiché $\mathbf{P} \mathbf{X} = \mathbf{0}$, allora $\mathbf{P} \mathbf{Y} = \mathbf{P} \mathbf{Y}_+$.

Definita \mathbf{P}_2 , è quindi possibile riscrivere la statistica LRT come

$$\text{LRT} = |\mathbf{Y}^T \mathbf{P} \mathbf{Y}| / |\mathbf{Y}^T \mathbf{P} \mathbf{Y}_+ \mathbf{Y}_+ \mathbf{P}_2 \mathbf{Y}_+| \quad (1.2)$$

ed essa ha distribuzione $A(p, n - q, g)$, sotto H_0 . Tale distribuzione si ottiene dal rapporto tra i determinanti di due matrice. Siano $\mathbf{A} \sim W_p(\mathbf{I}, m)$ e $\mathbf{B} \sim W_p(\mathbf{I}, n)$, indipendenti, dove $W_p(\mathbf{\Sigma}, m)$ indica una distribuzione di Wishart con matrice di scala $\mathbf{\Sigma}$ ed m gradi di libertà. Se $m \geq p$ allora

$$A = |\mathbf{A}| / |\mathbf{A} + \mathbf{B}| \sim A(p, m, n)$$

(si veda l'Appendice).

Per maggiori chiarimenti sulla distribuzione A e sulle distribuzioni ad essa collegate si rimanda al Mardia *et al.* (1979).

Nel caso in cui si voglia utilizzare questo test per dimostrare l'ipotesi nulla $H_0: \mathbf{C}_1 \mathbf{B} \mathbf{M}_1 = \mathbf{D}$, con $\mathbf{M}_1 \neq \mathbf{I}$, la statistica test che si ottiene è una generalizzazione del LRT e prende il nome di lambda di Wilks (si veda il Cap. 2).

1.2.2. UIT

L'ipotesi nulla $H_0: \mathbf{C}_1 \mathbf{B} \mathbf{M}_1 = \mathbf{D}$ è vera e se e solo se vale la relazione $\mathbf{b}^T \mathbf{C}_1 \mathbf{B} \mathbf{M}_1 \mathbf{a} = \mathbf{b}^T \mathbf{D} \mathbf{a}$, per ogni valore di \mathbf{a} e \mathbf{b} . Se si sostituiscono rispettivamente $\mathbf{b}^T \mathbf{C}_1$ e $\mathbf{M}_1 \mathbf{a}$ nelle relazioni

$$\mathbf{H} = \mathbf{M}_1^T \mathbf{Y}_+^T \mathbf{P}_2 \mathbf{Y}_+ \mathbf{M}_1 \quad \text{ed} \quad \mathbf{E} = \mathbf{M}_1^T \mathbf{Y}_+^T \mathbf{P} \mathbf{Y}_+ \mathbf{M}_1,$$

e scriviamo il rapporto tra \mathbf{H} ed \mathbf{E} , si ottiene

$$\frac{\{ \mathbf{b}^T \mathbf{C}_1 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}_+ \mathbf{M}_1 \mathbf{a} \}}{\{ \mathbf{b}^T \mathbf{C}_1 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}_1^T \mathbf{b} \} \{ \mathbf{a}^T \mathbf{M}_1^T \mathbf{Y}_+^T \mathbf{P} \mathbf{Y}_+ \mathbf{M}_1 \mathbf{a} \}} \quad (1.3)$$

che, sotto l'ipotesi nulla, si distribuisce come una $F_{1, n - q}$ moltiplicata per $(n - q)^{-1}$, per \mathbf{a} e \mathbf{b} fissati.

Massimizzando rispetto a \mathbf{b} si ottiene

$$\mathbf{a}^T \mathbf{H} \mathbf{a} / \mathbf{a}^T \mathbf{E} \mathbf{a},$$

che, sotto H_0 si distribuisce come una $[g/(n - q)]F_{g, n - q}$, con \mathbf{a} fissato (perché \mathbf{H} ed \mathbf{E} hanno una distribuzione di Wishart e sono tra loro indipendenti).

Infine, massimizzando rispetto ad \mathbf{a} si ottiene λ_1 , che è il più grande degli autovalori di $\mathbf{H}\mathbf{E}^{-1}$.

Posto $\theta = \lambda_1 / (1 + \lambda_1)$, θ rappresenta la più grande delle radici di $\mathbf{H}(\mathbf{H} + \mathbf{E})^{-1}$. L'ipotesi nulla viene rifiutata per valori elevati di θ , che segue, sotto H_0 , la distribuzione $\theta(r, n - q, g)$. Si ricorda che per $\mathbf{A} \sim W_p(\mathbf{I}, m)$ e $\mathbf{B} \sim W_p(\mathbf{I}, n)$, indipendenti con $m \geq p$, allora il più grande degli autovalori di $(\mathbf{A} + \mathbf{B})^{-1}\mathbf{B}$ ha distribuzione $\theta(m, n, p)$. (Si veda l'Appendice).

Nel caso in cui il rango della matrice \mathbf{M}_I sia $r = 1$, allora non è necessario massimizzare anche sotto \mathbf{a} , perché la statistica diventa semplicemente un rapporto tra scalari. In questo caso il test UIT diventa equivalente al test LRT.

Se invece il rango di \mathbf{C}_I è $g = 1$, allora la massimizzazione rispetto a \mathbf{b} non è necessaria, poiché in questo caso sia \mathbf{H} che $\mathbf{H}\mathbf{E}^{-1}$ hanno rango 1. Ancora una volta le due statistiche UIT e LRT sono equivalenti. Infatti, l'unico autovalore non nullo di $\mathbf{H}\mathbf{E}^{-1}$ è

$$\lambda_1 = \text{tr}(\mathbf{H}\mathbf{E}^{-1}) = \{\mathbf{C}_I(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}_I^T\}^{-1} \mathbf{C}_I(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}_+ \mathbf{M}_I \mathbf{E}^{-1} \mathbf{M}_I^T \mathbf{Y}_+^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}_I^T.$$

Inoltre, se fissiamo $\mathbf{d} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{C}_I^T$, allora sotto l'ipotesi nulla H_0 si ha che $\mathbf{M}_I^T \mathbf{Y}_+^T \mathbf{d} \sim N_p(\mathbf{0}, (\mathbf{d}^T \mathbf{d}) \mathbf{M}_I^T \boldsymbol{\Sigma} \mathbf{M}_I)$ e sotto entrambe le ipotesi, nulla e alternativa, che $\mathbf{E} \sim W_r(\mathbf{M}_I^T \boldsymbol{\Sigma} \mathbf{M}_I, n - q)$. Allora sotto l'ipotesi nulla, λ_1 ha distribuzione

$$\lambda_1 \sim (n - q)^{-1} T_{r, n - q}^2 = [r/(n - q - r + 1)] F_{r, n - q - r + 1},$$

dove T^2 indica la distribuzione di Hotelling (si veda l'Appendice).

Supponendo ora che \mathbf{B} sia il vero valore della matrice dei coefficienti, dal test UIT si ottiene che la probabilità che il rapporto descritto nella (1.2) sia minore o uguale di $\theta_\alpha / (1 - \theta_\alpha)$, per tutti i valori di \mathbf{a} e \mathbf{b} , è pari a $1 - \alpha$, con θ_α che indica il percentile di livello α della distribuzione $\theta(r, n - q, g)$.

1.3 Intervalli di confidenza simultanei per B

Utilizzando i risultati presentati nel paragrafo precedente è possibile ottenere degli intervalli di confidenza simultanei (SCIs) per le stime dei parametri \hat{B} , ossia tali per cui

$$P\left(b^T C_1 B M_1 a \in b^T C_1 (X^T X)^{-1} X^T Y M_1 a \pm \left\{ \frac{\theta_\alpha}{1 - \theta_\alpha} (a^T E a) \left[b^T C_1 (X^T X)^{-1} C_1^T b \right] \right\}^{1/2}, \forall a, b\right) = 1 - \alpha$$

Il valore $b^T C_1 (X^T X)^{-1} X^T Y M_1 a = b^T C_1 \hat{B} M_1 a$ è uno stimatore non distorto di $b^T C_1 B M_1 a$.

In alcune applicazioni a e/o b sono forniti a priori e i limiti dell'intervallo di confidenza possono essere ristretti. Ad esempio, se a è fissato, allora gli intervalli di confidenza per b si possono ottenere sostituendo $\theta_\alpha/(1 - \theta_\alpha)$ con $[g/(n-q)]F_{r, n-q-r+1; \alpha}$, che rappresenta il percentile di livello α di una F con r e $n-q-r+1$ gradi di libertà.

Se invece è b ad essere fissato, allora $\theta_\alpha/(1 - \theta_\alpha)$ può essere sostituito da

$$(n-q)^{-1} T_{r, n-q; \alpha}^2 = [r/(n-q-r+1)] F_{r, n-q-r+1; \alpha}.$$

Infine, se sono noti sia a che b si può ottenere un singolo intervallo di confidenza sostituendo $\theta_\alpha/(1 - \theta_\alpha)$ con $(n - q)^{-1} F_{1, n-q; \alpha}$.

1.4 Alcune considerazioni conclusive

1.4.1 Stima dei coefficienti nel caso in cui X non sia di rango pieno

Talvolta può succedere, per ragioni di simmetria nel disegno sperimentale, che la matrice X di dimensioni $(n \times q)$ delle esplicative non sia di rango pieno, ma abbia rango $k < q$. Ovviamente, in questa situazione, la matrice $X^T X$ non è più invertibile e quindi non è possibile calcolare le stime dei coefficienti di regressione. In questo caso, per risolvere il problema, può essere utile considerare una partizione della matrice X , in modo da restringere la regressione ad un gruppo di colonne della matrice delle esplicative. In particolare, si considera $X = (X_1, X_2)$, dove X_1 è la matrice di dimensioni $(n \times k)$ di rango pieno. Allora X_2 può essere scritta come $X_2 = X_1 A$, dove A è una generica matrice di dimensione $(k \times (q-k))$. Analogamente, anche per la matrice B si considera la partizione $B^T = (B_1^T, B_2^T)$.

Il modello può essere dunque riformulato nel seguente modo

$$Y = X_I B^*_I + U_I,$$

dove B^*_I indica la somma $B_I + AB_2$ univocamente determinata. Effettuando queste sostituzioni, è possibile ottenere comunque una stima dei coefficienti, nonostante la matrice delle esplicative X non sia di rango pieno, poiché il modello $Y = X_I B^*_I + U_I$ permette di utilizzare gli usuali metodi di stima per B^*_I , che rappresenta una riparametrizzazione di B .

1.4.2 Valutazione dell'effetto della media sulle risposte

Nel modello considerato talvolta può essere utile separare l'effetto della media μ sulle risposte, dall'effetto delle altre variabili esplicative indipendenti. Il modello può essere riscritto come

$$Y = \mathbf{1}\mu^T + XB + U, \quad (1.4)$$

dove $(\mathbf{1}, X)$ rappresenta una matrice $(n \times (1 + q))$. Allora

$$\left[\begin{pmatrix} \mathbf{1}^T \\ X^T \end{pmatrix} (\mathbf{1} \quad X) \right]^{-1} = \begin{pmatrix} n^{-1} & 0 \\ 0 & (X^T X)^{-1} \end{pmatrix},$$

poiché $X^T \mathbf{1} = 0$. Allora, poiché vale la (1.1), si ha $\bar{\mu} = \bar{y}$, con \bar{y} vettore di medie calcolate su tutti i dati, che è indipendente dalla stima di B che si ottiene come

$$\hat{B} = (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T (Y - \mathbf{1}\bar{y}^T).$$

In particolare, se la media di Y è stimata dal vettore \bar{y} , allora la (1.4) si può anche scrivere come

$$Y - \mathbf{1}\bar{y}^T = XB + Z,$$

dove con Z , si intende $Z = U - \mathbf{1}\bar{u}^T$, che rappresenta la matrice degli errori centrata. Ogni vettore colonna di tale matrice ha distribuzione $N_{np}(\mathbf{0}, \Sigma \otimes Q)$, dove $Q = \mathbf{I} - n^{-1} \mathbf{1}\mathbf{1}^T$. Con il simbolo \otimes si indica il prodotto di Kronecker. Sia A una matrice $(m \times n)$, con elementi a_{ij} e B una matrice $(p \times q)$ di elementi b_{kl} , allora il prodotto di Kronecker di A e B è definito come

$$A \otimes B = \begin{bmatrix} a_{11}B & a_{12}B & \cdot & \cdot & \cdot & a_{1n}B \\ a_{12}B & a_{22}B & \cdot & \cdot & \cdot & a_{2n}B \\ \cdot & \cdot & & & & \cdot \\ \cdot & \cdot & & & & \cdot \\ \cdot & \cdot & & & & \cdot \\ a_{m1}B & a_{m2}B & \cdot & \cdot & \cdot & a_{mn}B \end{bmatrix}$$

che è una matrice ($mp \times nq$).

Quanto descritto rappresenta un'estensione di quanto detto in precedenza, nel caso in cui, nello stimare i parametri B , si voglia valutare separatamente l'effetto della media μ , rispetto alle altre esplicative e quindi capire qual è il peso di μ nella procedura di determinazione di \hat{B} .

1.4.3 Dalle stime OLS alle stime GLS

Nel modello univariato, il calcolo delle stime dei coefficienti di regressione si possono anche basare sul metodo dei minimi quadrati ordinari, sotto l'ipotesi che i termini d'errore abbiano media nulla e siano omoschedastici (ipotesi di Gauss-Markov). Quando la matrice di varianza e covarianza degli errori non può più essere supposta uguale a $\sigma^2 I$, lo stimatore OLS non è più BLUE (*Best Linear Unbiased Estimator*). In queste situazioni si utilizza la stima GLS (*Generalized Least Squares*), che si basa su una trasformazione della matrice di covarianza. Partendo dalle assunzioni

$$E(u) = 0 \text{ e } V(u) = \Omega,$$

con Ω supposta nota, si considera la trasformazione del modello

$$z = \Omega^{-1/2} X B + v, \tag{1.5}$$

dove

$$z = \Omega^{-1/2} y \quad \text{e} \quad v = \Omega^{-1/2} u.$$

Il modello (1.5) soddisfa le ipotesi di Gauss-Markov e quindi permette di calcolare lo stimatore non distorto di B .

Nel caso multivariato, ciò che si verifica è che la stima OLS, che si ottiene se $\Sigma = \sigma^2 I$ e la stima GLS conducono allo stesso stimatore. Riscrivendo il modello in forma vettoriale, si ottiene infatti

$$\mathbf{Y}^V = \mathbf{X}^* \mathbf{B}^V + \mathbf{U}^V,$$

dove $\mathbf{Y}^V = (\mathbf{y}_{(1)}^T, \dots, \mathbf{y}_{(p)}^T)$ e $\mathbf{X}^* = \mathbf{I}_p \otimes \mathbf{X}$. Per quanto riguarda il termine d'errore \mathbf{U}^V , si assume che esso abbia media $\mathbf{0}$ e matrice di varianza e covarianza

$$\mathbf{V}(\mathbf{U}^V) = \mathbf{\Omega} = \mathbf{\Sigma} \otimes \mathbf{I}_p.$$

Allora lo stimatore GLS di \mathbf{B} si ottiene come

$$\begin{aligned} \hat{\mathbf{B}}^V &= (\mathbf{X}^{*\top} \mathbf{\Omega}^{-1} \mathbf{X}^*)^{-1} \mathbf{X}^{*\top} \mathbf{\Omega}^{-1} \mathbf{Y}^V \\ &= (\mathbf{\Sigma} \otimes \mathbf{X}^T \mathbf{X})^{-1} (\mathbf{\Sigma}^{-1} \otimes \mathbf{X}^T) \mathbf{Y}^V \\ &= [\mathbf{I} \otimes (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{Y}^V. \end{aligned}$$

In particolare, la stima OLS si ottiene nel caso in cui $\mathbf{\Sigma} = \mathbf{I}$; in questa situazione lo stimatore OLS e lo stimatore GLS coincidono.

È possibile valutare il comportamento dello stimatore GLS asintoticamente, ossia nell'ipotesi in cui n , che rappresenta la numerosità campionaria, tenda ad ∞ .

Sia $\mathbf{x}_1, \mathbf{x}_2, \dots$ una sequenza di variabili indipendenti fissate e sia u_1, u_2, \dots una sequenza di termini d'errore casuali. Siano $\mathbf{X} = \mathbf{X}_n = (\mathbf{x}_1^T, \mathbf{x}_2^T, \dots, \mathbf{x}_n^T)$, $\mathbf{\Omega} = \mathbf{\Omega}_n$ e $\mathbf{B} = \mathbf{B}_n$. Si supponga inoltre che

$$\lim_{n \rightarrow \infty} (\mathbf{X}_n^T \mathbf{\Omega}_n^{-1} \mathbf{X}_n)^{-1} = \mathbf{0}.$$

Allora il valore atteso della stima di \mathbf{B} sarà \mathbf{B} per ogni valore di n considerato. Inoltre la varianza dello stimatore tenderà a 0 quando $n \rightarrow \infty$.

Ciò implica che lo stimatore di \mathbf{B} è consistente.

Capitolo 2

MANOVA: Analisi della varianza multivariata

Nel caso della regressione multipla, una delle procedure più utilizzate per valutare l'uguaglianza delle medie in gruppi differenti e per valutare la significatività dell'effetto di un'esplicativa sulla variabile risposta è l'analisi della varianza (**ANOVA**).

Anche nel caso multivariato si può utilizzare una procedura simile, la MANOVA (*multivariate analysis of variance*). Ad esempio, nel caso in cui si abbiano due variabili dipendenti e si voglia testare se entrambe sono influenzate da una certa esplicativa è possibile calcolare una statistica F , non più univariata. Ciò che si ottiene è una statistica F multivariata (lambda di Wilks) basata sul confronto tra la matrice di varianza e covarianza degli errori e la matrice di varianza e covarianza delle esplicative. Infatti, la formula della statistica F non è basata solo sulla somma dei quadrati tra ed entro i gruppi, ma anche dai prodotti incrociati (*crossproducts*).

Il calcolo della statistica F permette di rispondere alla domanda "Il modello è significativo per ogni variabile dipendente?". La covarianza deve essere inclusa perché le due variabili dipendenti potrebbero essere correlate ed è necessario testare se tale correlazione è statisticamente significativa. Se le due variabili risposta sono fortemente correlate si ottengono informazioni ridondanti e ciò è espresso dalla covarianza.

L'ipotesi nulla che si andrà a testare riguarda l'assenza di differenze in media in ogni variabile dipendente per gruppi differenti formati a seconda dei livelli e delle categorie delle esplicative.

Oltre a valutare se le risposte sono influenzate contemporaneamente dalle esplicative, può essere interessante valutare se le variabili dipendenti prese singolarmente vengono influenzate dalle indipendenti considerate. Infatti, oltre a calcolare la statistica F multivariata che considera gli effetti principali e le interazioni, è possibile esaminare i test F univariati per ogni singola variabile ed interpretarli singolarmente.

I test considerati nel caso univariato si possono estendere al caso multivariato. Se in precedenza si aveva il rapporto tra devianza di regressione e devianza residua, nel caso multivariato si moltiplica l'inverso della matrice E per la matrice H , dove la matrice E , matrice di devianza residua, si ottiene come

$$E = \mathbf{u}_l^T \mathbf{u}_l = (n_1 - 1) S_1^2 + (n_2 - 1) S_2^2 + \dots + (n_p - 1) S_p^2$$

con S_1^2, S_2^2, S_p^2 matrici di varianza e covarianza all'interno dei rispettivi gruppi, n_1, n_2, \dots, n_p numerosità dei gruppi, e H , matrice di devianza tra i gruppi, data da

$$H = \sum_j n_j \beta_{j0}^T \mathbf{B}_{jl}^T \mathbf{B}_{jl} \beta_{j0} = [\boldsymbol{\beta}_l - \mathbf{B}_l]^T [\boldsymbol{\beta}_l - \mathbf{B}_l] [\mathbf{y}_l - \mathbf{m}_l]^T [\mathbf{y}_l - \mathbf{m}_l].$$

La somma delle due matrici, in genere indicata con $T = E + H$, è la matrice di devianza totale.

I test, che di seguito verranno descritti, permettono di verificare la presenza di una relazione di dipendenza lineare della matrice delle risposte Y dalla matrice delle esplicative X .

2.1 Test multivariati

Esistono alcuni test multivariati che permettono di focalizzare l'attenzione sulle variabili indipendenti e sulle loro interazioni. In particolare ci sono quattro tipi di test utilizzati per questo scopo.

2.1.1 T^2 di Hotelling

Il test T^2 di Hotelling è definito come

$$T^2 = \text{tr}(E^{-1} H) = \sum \lambda_{ij},$$

dove λ indica l'autovalore di una matrice \mathbf{A} ($k \times k$). È il test più utilizzato e tradizionale basato su due gruppi formati dalle variabili indipendenti. Ad esso è correlata la statistica che prende il nome di *Hotelling's Trace*. Per trasformare il coefficiente *Trace* nel coefficiente T^2 si moltiplica il *Trace* per $(n - g)$, dove n è la numerosità complessiva e g il numero dei gruppi. Il T^2 avrà lo stesso valore della statistica F , lo stesso numero di gradi di libertà e lo stesso livello di significatività del coefficiente *Trace*.

2.1.2 Lambda di Wilks

La lambda di Wilks si può scrivere come

$$\Lambda = |\mathbf{E}| / |\mathbf{H} + \mathbf{E}| = \prod (1 / (1 + \lambda_{ij})) \quad (2.1)$$

oppure, analogamente, come

$$\Lambda = \frac{|\mathbf{E}|}{|\mathbf{H} + \mathbf{E}|} = \frac{\sum_{l=1}^p \sum_{j=1}^q (\beta_{jl} \mathbf{B}_{jl} - \boldsymbol{\mu})(\beta_{jl} \mathbf{B}_{jl} - \boldsymbol{\mu})^T}{\sum_{l=1}^p \sum_{j=1}^q (\mathbf{y}_{lj} - \boldsymbol{\mu})(\mathbf{y}_{lj} - \boldsymbol{\mu})^T},$$

per il quale ci sono più di due gruppi di variabili indipendenti. È una misura delle differenze tra i gruppi del vettore di medie sulle variabili indipendenti. Più piccolo è il valore della statistica Λ , più grandi saranno le differenze esistenti. Per calcolare la significatività di Λ viene spesso utilizzata una sua trasformazione, la V di Bartlett. Usate in combinazione, le due statistiche, la Λ di Wilks e la V di Bartlett, consentono il calcolo di un test per le differenze in media nella MANOVA, nel caso in cui ci siano più variabili dipendenti e più di due gruppi formati dalle esplicative. I test considerati precedentemente (t -test, Hotelling's T e test F) sono tutti casi particolari della Λ di Wilks. In particolare, in casi particolari Λ segue una distribuzione F , ossia quando

- $q = 2$ e $l \geq 2$, allora Λ si distribuirà come una F con $l - 2$ e $2 \times (n - q - 1)$ gradi di libertà;
- se $q \geq 1$ e $l = 2$, allora Λ si distribuirà come una F con p e $n - q - 1$ gradi di libertà.

Più in generale sotto l'ipotesi nulla H_0 e con un'elevata numerosità campionaria, la statistica test si distribuisce come un χ^2 con $p \times (l - 1)$ gradi di libertà.

2.1.3 Pillai-Bartlett trace

Il test Pillai-Bartlett trace è definito come

$$\mathbf{PB} = \text{Tr}(\mathbf{E}(\mathbf{E} + \mathbf{H})^{-1}),$$

che si incontra anche nella MDA (*multiple discriminant analysis*), ossia in quella parte della MANOVA, simile all'Analisi delle Componenti Principali, che ha come obiettivo la classificazione. Tale test è il più robusto di tutti, ossia il più affidabile anche nel caso in cui certe ipotesi siano violate.

2.1.4 GCR

Il test GCR (*Roy's greatest characteristic root*) è definito come

$$\mathbf{GCR} = \max \lambda_{ij}. \quad (2.2)$$

Esso è simile al Pillai-Bartlett, ma è basato sulla prima e più importante radice. Questo test è meno robusto dei precedenti a causa della violazione delle assunzioni di normalità multivariata.

2.1.5 Alcune conclusioni

I test presentati in questo paragrafo, permettono di riportare il modello multivariato al caso univariato in termini di verifica d'ipotesi. In particolare, essi permettono di verificare se le variabili dipendenti \mathbf{Y} dipendono simultaneamente dalle esplicative \mathbf{X} . Infatti, l'ipotesi alternativa è quella dell'esistenza di una relazione lineare significativa tra \mathbf{Y} e \mathbf{X} . Si dimostra che \mathbf{H} è stima distorta della devianza fra i gruppi, mentre \mathbf{E} è stima corretta delle devianze fra i gruppi residui.

Nel caso più semplice in cui abbiamo solo due gruppi da confrontare, estendiamo l'ipotesi nulla del modello univariato

$$H_0: \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$$

al caso multivariato, in cui vogliamo verificare che i vettori di medie siano uguali, ossia

$$H_0 : \begin{pmatrix} \mu_{11} \\ \mu_{21} \\ \vdots \\ \mu_{p1} \end{pmatrix} = \begin{pmatrix} \mu_{12} \\ \mu_{22} \\ \vdots \\ \mu_{p2} \end{pmatrix}.$$

L'estensione al caso multivariato del test t è data da

$$T^2 = \frac{n_1 n_2}{n_1 + n_2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2),$$

con

$$T^2 \sim F_{n_1+n_2-p-1, (n_1+n_2-2)p},$$

dove \mathbf{S} è la matrice di covarianza “unita” (*pooled*) entro i gruppi e n_1 e n_2 sono le numerosità dei due gruppi.

Nel caso più generico in cui abbiamo p gruppi di osservazioni possiamo testare l'ipotesi nulla $H_0: \tau_1 = \tau_2 = \dots = \tau_p = 0$, dove τ_i indica il vettore di medie dell' i -esimo gruppo. In questa situazione si utilizza il Λ di Wilks (vedi formula **(2.1)**).

Si possono utilizzare due tipi di approcci per valutare la significatività dei test effettuati e per confrontare la devianza dei trattamenti e la devianza degli errori o residua. Il primo metodo è l'approccio della verosimiglianza; il secondo è quello delle combinazioni lineari.

Il metodo della verosimiglianza si serve del Λ di Wilks e rifiuta l'ipotesi nulla H_0 se il valore della statistica test è grande.

Il secondo approccio, invece, si basa sul test GCR di Roy (vedi formula **(2.2)**). Tale statistica tenta di trovare combinazioni lineari delle osservazioni che massimizzino il rapporto tra varianza spiegata e varianza residua, massimizzando la varianza tra i gruppi e contemporaneamente minimizzando quella entro i gruppi.

Così come l'ANOVA testa le differenze in media della variabile dipendente per varie categorie delle variabili indipendenti, la MANOVA testa la differenze nei vettori di medie delle risposte per le diverse categorie delle esplicative.

Inoltre, così come nell'ANOVA la varianza spiegata è la stima distorta sotto l'ipotesi alternativa H_1 della quantità $(p - 1) \sigma^2$, essendo corretta sotto l'ipotesi nulla H_0 , ora nel caso della MANOVA la matrice \mathbf{B} è la stima corretta della matrice di devianza tra i gruppi solo sotto l'ipotesi nulla H_0 .

2.2 Test “post-hoc”

Il passo successivo al calcolo dei test nella procedura di analisi della varianza multivariata è rappresentato dai test a posteriori (*post-hoc*) che si eseguono nel caso in cui il test F mostri che il vettore di medie non sia lo stesso nei diversi gruppi formati sulla base della categorie delle esplicative. Per effettuare il test *post-hoc* si calcolano le statistiche F univariate che sono usate per determinare quali medie dei gruppi differiscono significativamente dalle altre. Allo stesso tempo si eseguono test (confronti multipli) su coppie di gruppi per verificarne somiglianze e differenze.

I test utilizzati per i confronti a posteriori sono:

- il test di **Bonferroni**, utilizzato se ci sono pochi gruppi, poiché produce le differenze in media delle variabili dipendenti tra coppie di gruppi (ad esempio differenze nei punteggi ottenuti nei test per ogni coppia di gruppi che differiscono per scuola di provenienza);
- il test di **Tukey**, preferibile se i gruppi sono numerosi;
- altri test utilizzati quando le assunzioni di omogeneità della varianza sono violate.

Queste procedure rappresentano le fasi conclusive della MANOVA e permettono di capire meglio le differenze rilevate tra i gruppi.

2.3 Assunzioni alla base dell'analisi multipla della varianza

La procedura di analisi multipla della varianza si basa su una serie di assunzioni che, se violate, rendono non affidabili gli indicatori calcolati.

Tali assunzioni sono:

- l'indipendenza delle osservazioni, poiché nell'ipotesi in cui questa condizione non valga, allora la MANOVA non è robusta;
- dimensioni uguali per tutti i gruppi;
- somme dei quadrati appropriate;
- dimensioni adeguate del campione;
- distribuzione casuale dei residui;
- omoschedasticità (omogeneità delle varianze e delle covarianze): la varianza di ogni variabile dipendente continua deve essere simile, come testato dal test di Levene (si veda Levene, 1960 [4]). Tale test, simile a quello di Bartlett, è utilizzato su k campioni per verificare se siano omogenei in varianza, ossia se hanno la stessa varianza. Viene utilizzato quindi per testare ipotesi del tipo

$$H_0: \sigma_1 = \sigma_2 = \dots = \sigma_p \quad \text{vs} \quad H_1: \sigma_i \neq \sigma_j \quad \text{per almeno un } i \neq j.$$

Se si hanno n osservazioni per una variabile Y , suddivisa in k gruppi di numerosità n_i , la statistica di Levene prende la seguente forma

$$W = \frac{(n-k) \sum_{i=1}^k n_i (\bar{Z}_i - \bar{Z})^2}{(k-1) \sum_{i=1}^k \sum_{j=1}^{n_i} (Z_{ij} - \bar{Z}_i)^2},$$

con $Z_{ij} = |Y_{ij} - \bar{Y}_i|$, dove \bar{Y}_i è la media dell' i -esimo gruppo.

Per ognuno dei k gruppi formati dalle variabili indipendenti, la covarianza tra le coppie di variabili dipendenti deve essere la stessa. Per verificare tale assunzione si può calcolare il test M di Box (si veda Anderson, 1958 [1]; Box, 1949 [2]; Seber, 1984 [8]) che verifica l'assunzione di omoschedasticità utilizzando un test F . Tale test è basato sul rapporto di verosimiglianza e assume la forma

$$M = (n-g) \log |S| - \sum_{i=1}^g (n_i-1) \log |S_i|,$$

dove S è la matrice di varianza "pooled", n è la numerosità complessiva, n_i è la numerosità dell' i -esimo gruppo, g è il numero di gruppi e S_i è la matrice di varianza e covarianza dell' i -esimo gruppo.

Se il p -value associato al test è inferiore a 0.05 allora le covarianze sono significativamente diverse;

- omogeneità della regressione: i coefficienti delle covariate sono gli stessi per ogni gruppo formato da variabili categoriali e misurato sulle variabili dipendenti;
- distribuzione normale multivariata: si assume la normalità multivariata se ogni variabile segue la distribuzione normale. La MANOVA è robusta anche in caso di violazione di questa assunzione se la dimensione del campione non è piccola. Per verificare tale assunzione si possono utilizzare procedure diverse: è possibile rappresentare graficamente i dati e verificare la presenza di simmetria, calcolare la curtosi, valutare la presenza di valori anomali; si può utilizzare il test di Shapiro-Wilks che considera come ipotesi nulla la normalità dei dati; si può rappresentare il diagramma "quantile-quantile", che confronta i quantili della distribuzione normale con quelli dei dati. Queste procedure sono quelle utilizzate anche nel caso univariato perché si assume che se tutte le variabili sono normali allora complessivamente avremo una distribuzione normale multivariata.

Altri controlli si possono effettuare calcolando i quadrati delle distanze dalla media delle osservazioni poiché se la popolazione è normale e n e $n - p$ sono maggiori di 30, allora tali distanze dovrebbero seguire la distribuzione di un

$$\chi_{n-p}^2.$$

- assenza di valori anomali;
- le covariate sono linearmente correlate o in una relazione nota con le variabili dipendenti, la forma della relazione tra covariate e variabili risposta deve essere nota. Spesso le covariate sono trasformate per stabilire una relazione lineare.

Per ulteriori chiarimenti sulla MANOVA e sui test ad essa correlati si rimanda a Miles, 2003 [5]; Vittadini, 1999 [9]; Schatz, 2006 [7].

2.4 Utilizzi della MANOVA

La MANOVA può avere molteplici utilizzi. La procedura può essere utilizzata per confrontare gruppi formati da variabili indipendenti di tipo categoriale rispetto alle differenze nei gruppi su un insieme di variabili risposta continue; per usare l'assenza di differenze per un gruppo di variabili dipendenti come un criterio per rendere più piccolo e facilmente gestibile il numero di variabili; per identificare le variabili indipendenti che differenziano maggiormente le risposte.

Nell'ipotesi in cui non ci sia correlazione tra le variabili risposta, è accettabile l'utilizzo di una serie di test basati sull'ANOVA univariata piuttosto che l'utilizzo della MANOVA. Nella maggior parte dei casi, però, le variabili sono tra loro correlate ed in questo caso è preferibile l'utilizzo della MANOVA, poiché tale procedura è sensibile non solo alle differenze in media, ma anche al segno e alla forza della correlazione esistente tra le variabili dipendenti.

Nella procedura della MANOVA non c'è un limite al numero di variabili dipendenti che possono essere inserite, ma più cresce il numero di risposte e più diventa difficile l'interpretazione dei risultati ottenuti. Inoltre, si riduce la potenza dei test poiché cresce la probabilità di commettere errori del II Tipo.

2.5 Conclusioni

La procedura di analisi della varianza multivariata permette di estendere al caso di regressione multivariata le conclusioni a cui si giungeva nel caso univariato e di applicare ai dati una serie di test per eseguire verifiche d'ipotesi.

La MANOVA può essere eseguita da una serie di pacchetti statistici. Ad esempio in SAS si usa il comando `Mtest`; in SPSS è presente l'opzione MANOVA.

Anche in **R** sono stati sviluppati alcuni pacchetti che si utilizzano per la regressione multivariata e che danno in *output* i test citati con i rispettivi *p-values*.

Tali pacchetti saranno presentati nel prossimo capitolo, assieme all'analisi di un insieme di dati reali.

Capitolo 3

Un caso studio

3.1 I dati

Il modello di analisi multivariata **(1)** può essere applicato in diversi campi nel caso in cui si voglia valutare l'effetto di alcune variabili esplicative X su più variabili risposta Y .

In particolare l'applicazione che si prende in considerazione in questo capitolo riguarda dati di natura antropometrica raccolti su $n = 3028$ unità statistiche. Tali unità sono rappresentate da adulti ed anziani sardi.

Le analisi riportate di seguito sono state eseguite con il software **R**, (www.r-project.org) che consente di eseguire analisi grafiche e calcoli statistici grazie ad alcune funzioni contenute all'interno di sue librerie. Tramite i pacchetti di base e quelli disponibili gratuitamente in rete, **R** consente di stimare modelli lineari e modelli lineari generalizzati, di eseguire regressioni non lineari, di analizzare serie storiche e di effettuare test parametrici e non-parametrici. Il *software* consente anche di costruire funzioni esterne e semplici programmi, di effettuare regressioni robuste, multivariate e non parametriche, test statistici classici e controllo della qualità.

Nel dataset analizzato, l'interesse degli antropologi è particolarmente rivolto alla popolazione con età superiore ai 65 anni, per la quale si vuole valutare, separatamente per i due sessi, l'evoluzione che le variabili antropometriche hanno al crescere dell'età.

In particolare, si vuole valutare l'effetto su 23 variabili risposta di età, sesso e patologia, dove *patologia* è una variabile che assume valore 0 se l'unità statistica è sana, 1 se è diabetica e 2 se soffre di Alzheimer.

La variabile *sesso* è codificata come un fattore a due livelli che assume valore 1 per gli uomini e 2 per le donne.

La variabile *eta* misura l'età in anni delle unità statistiche al momento della rilevazione.

Le 23 variabili risposta considerate sono:

- *peso*, che rappresenta il peso (*kg*);
- *statura*, che misura l'altezza degli individui (*cm*);
- *pm.vita*, *pm.fianchi* e *pm.polpacc*, che misurano, rispettivamente, il perimetro della vita, dei fianchi e del polpaccio (*cm*);
- *pm.Br.contr*, che misura il perimetro del braccio contratto (*cm*);
- *dm.biacrom*, *dm.omero* e *dm.bicres* che misurano, rispettivamente, diametro biacromiale, dell'omero e bicrestiliaco (*cm*);
- *pl.Bicip*, *pl.tric*, *pl.sottosc* e *pl.soprail*, che indicano plica al bicipite, al tricipite, sottoscapolare e soprailiaca e sono misurate in *cm*;
- *lung.cefal*, che misura la lunghezza cefalica (*cm*);
- *R/H*, indicatore della resistenza normalizzata per la statura;
- *X_C/H*, indicatore della reattanza normalizzata per la statura, dove la reattanza rappresenta la forza che un condensatore oppone al passaggio di una corrente elettrica;
- *FFMLohman*, che misura la massa priva di grassi;
- *FMLohman*, che indica la massa grassa;
- *BMI (Body Mass Index)*, l'indice di massa corporea che si calcola come $\text{peso}/\text{statura}^2$;
- *WHR*, rapporto tra il perimetro della vita e quello dei fianchi;
- *STSR*, rapporto tra plica sottoscapolare e plica del tricipite;
- *AMA*, area muscolare del braccio, che si ottiene da una rielaborazione tra le misure del perimetro del braccio e la plica del tricipite;

- Angolo di fase, che si ottiene dall'arcotangente del rapporto tra X_c e R , trasformata in gradi.

Prima di effettuare qualunque tipo di analisi è necessario fattorizzare le due variabili

sexo e patologia:

```
dati$patologia<-factor(dati$patologia)
```

```
dati$sexo<-factor(dati$sexo)
```

```
contrasts(sexo)
```

```
  2
```

```
1 0
```

```
2 1
```

```
contrasts(patologia)
```

```
  1 2
```

```
0 0 0
```

```
1 1 0
```

```
2 0 1
```

Per effettuare le analisi di interesse è opportuno selezionare dal dataset le unità statistiche con età superiore a 65 anni.

Il nuovo dataset contiene, quindi, 743 osservazioni.

Sui dati considerati si possono calcolare le statistiche descrittive per le variabili considerate. Delle unità statistiche del dataset, il 72,7% non presenta patologie, il 19,9% soffre di diabete e il 7,5% di Alzheimer. Inoltre, gli individui di sesso maschile sono il 47% del totale, mentre il 53% sono di sesso femminile. Nel seguito sono riportate alcune statistiche descrittive di base per le variabili risposta.

	<u>eta</u>	<u>peso</u>	<u>statura</u>	<u>pm.vita</u>	<u>pm.fianchi</u>	<u>pm.polpac</u>
Minimo	65.04	30.00	124.6	50.00	63.1	21.10
1 st Qu	70.98	56.75	146.7	87.00	96.0	31.30
Mediana	76.82	67.00	153.0	95.50	102.2	34.00
Media	77.93	66.73	153.4	95.22	102.9	33.79
3 rd Qu.	83.93	76.10	160.0	102.20	108.6	36.40
Std Dev	8.21	14.45	9.58	150.00	148.0	55.00
Massimo	101.38	124.00	178.4	12.24	10.30	3.71
NA's				2	5	22

	<u>pm.br.contr</u>	<u>dm.biacrom</u>	<u>dm.omero</u>	<u>dm.bicres</u>	<u>pl.bicip</u>
Minimo	17.80	20.40	4.60	20.10	2.00
1 st Qu	26.50	31.00	6.10	28.90	8.00
Mediana	29.30	35.20	6.60	30.40	12.00
Media	29.32	34.20	6.65	30.23	13.97
3 rd Qu	32.00	37.83	7.20	32.10	18.00
Std Dev	4.02	3.69	0.69	2.36	7.99
Massimo	42.30	45.00	9.00	39.00	52.00
NA's	21	275	22	275	1

	<u>pl.tric</u>	<u>pl.sottosc</u>	<u>pl.soprail</u>	<u>lung.cef</u>	<u>rh</u>	<u>xch</u>
Minimo	2.60	2.00	2.00	16.70	188.4	14.43
1 st Qu	14.00	16.00	18.00	18.50	281.9	29.14
Mediana	21.00	21.50	25.00	19.00	323.8	34.42
Media	22.04	22.87	26.06	19.01	331.3	35.63
3 rd Qu	28.75	28.38	32.00	19.50	372.2	40.86
Std Dev	10.16	9.52	10.52	0.47	62.06	8.42
Massimo	70.00	62.00	70.00	21.00	571.4	75.38
NA's	2	5	21	477	90	90

	<u>ffmlohman</u>	<u>fmlohman</u>	<u>bmi</u>	<u>whr</u>	<u>stsr</u>	<u>ama</u>	<u>angolo_fase</u>
Minimo	24.99	4.612	16.20	0.57	0.27	7.13	3.17
1 st Qu	36.75	17.95	24.72	0.88	0.81	25.72	5.34
Mediana	41.77	23.251	27.94	0.93	1.05	33.42	6.13
Media	42.68	24.018	28.20	0.93	1.16	34.63	6.19
3 rd Qu	48.20	29.13	31.27	0.98	1.39	42.75	7.01
Std Dev	7.77	8.06	4.84	0.08	0.52	12.22	1.19
Massimo	70.31	56.31	45.18	1.29	4.18	84.64	12.42
NA's	90	90	4	6	10	26	90

Alcune delle variabili del dataset presentano un numero elevato di valori mancanti (NA's). In questa situazione si può decidere di sostituire i valori mancanti con la mediana o con la media della distribuzione, a seconda che ci sia o meno simmetria, oppure di eliminare la variabile. Nel caso di variabili come `lung.cefal`, `dm.biacrom` e `dm.bicres`, sembra essere più opportuno escludere le variabili dall'analisi poiché la prima ha il 67% dei dati mancanti, mentre le ultime due il 37%.

Per quanto riguarda le altre variabili, si è deciso di sostituire i valori mancanti con la mediana della distribuzione.

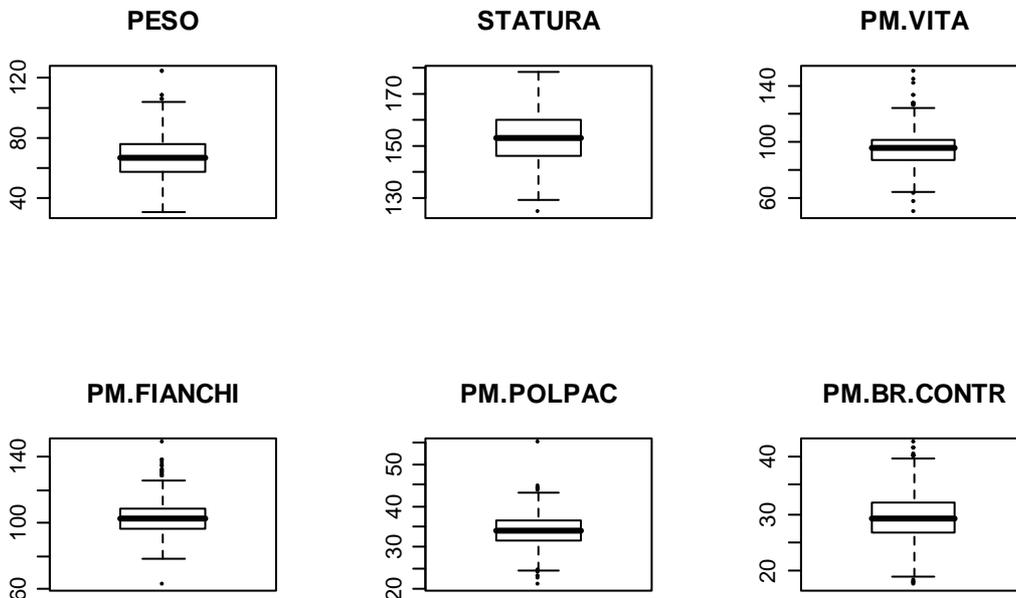
Il dataset studiato contiene pertanto 20 variabili risposta (in seguito all'eliminazione di `dm.bicres`, `dm.biacrom` e `lung.cef`) che devono essere valutate in base a sesso, età e patologia. Una prima analisi grafica permette di valutare se esiste una dipendenza effettiva delle variabili dipendenti dalle esplicative. Inoltre, l'analisi grafica consente di valutare se le variabili siano o meno simmetriche e se si possono ritenere normali. Utilizzando alcune librerie di **R** è possibile valutare, anche visivamente, la forza delle correlazioni esistenti tra le variabili risposta.

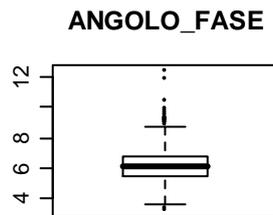
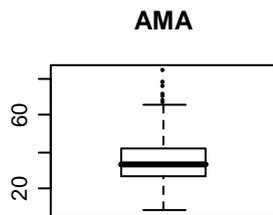
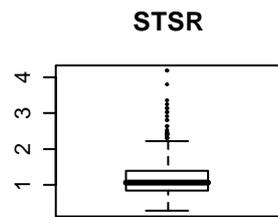
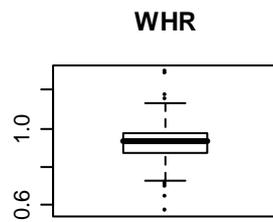
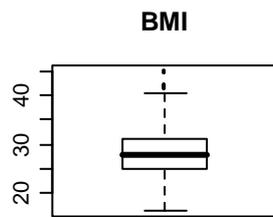
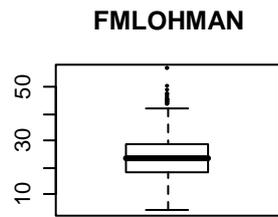
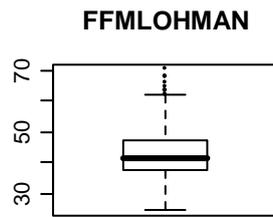
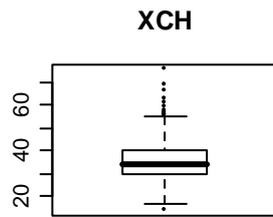
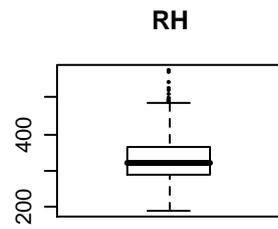
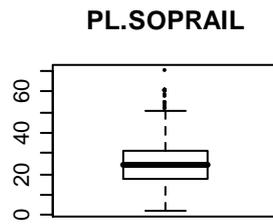
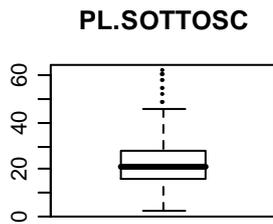
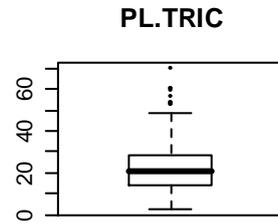
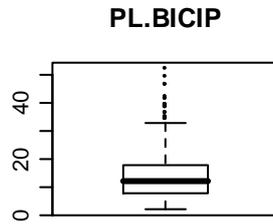
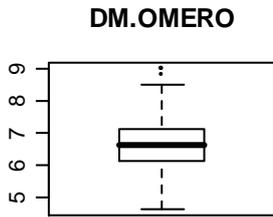
3.2 Analisi grafiche univariate

Per iniziare l'analisi sulla distribuzione delle variabili si possono considerare i boxplot, che permettono di visualizzare minimo, massimo, primo e terzo quartile e mediana della distribuzione.

Il comando di **R** per costruire boxplot è il seguente:

```
boxplot(x, data, subset, main="...")
```



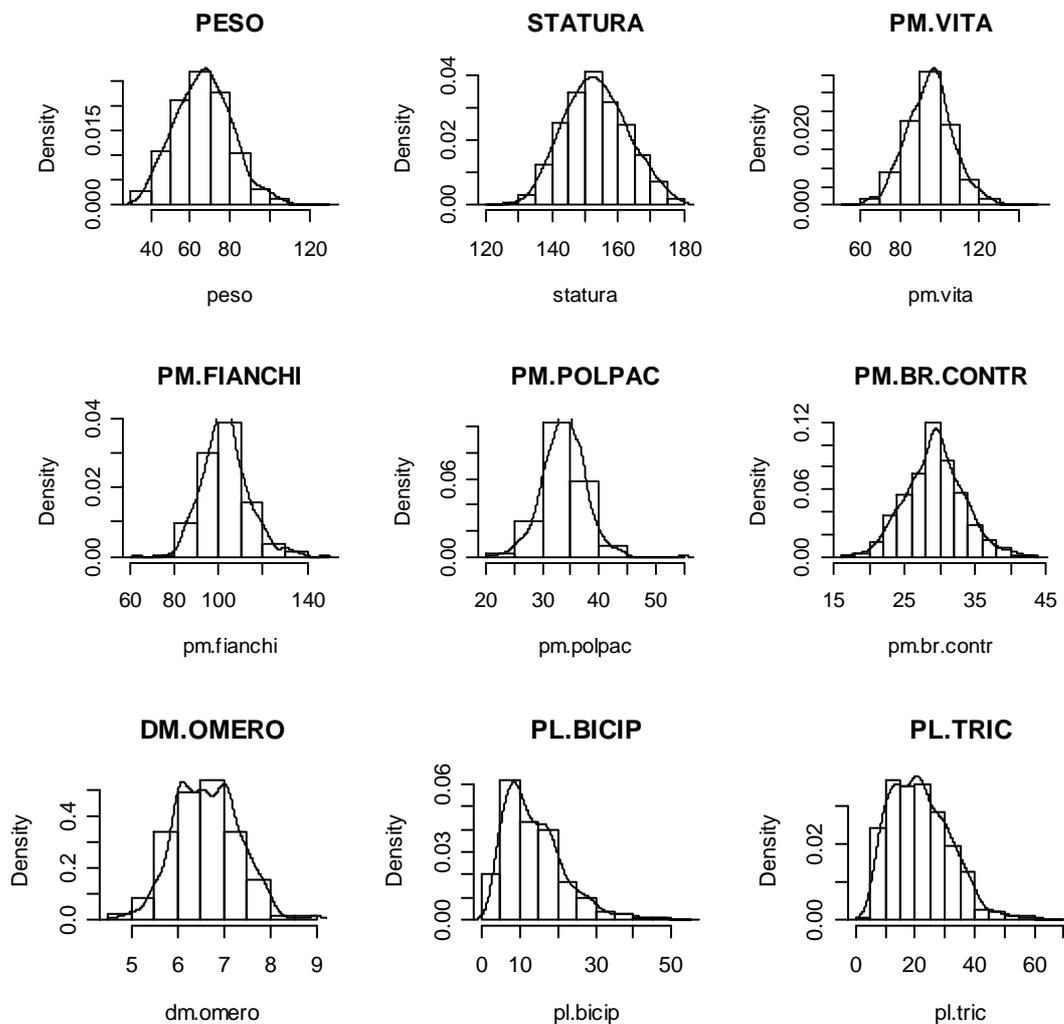


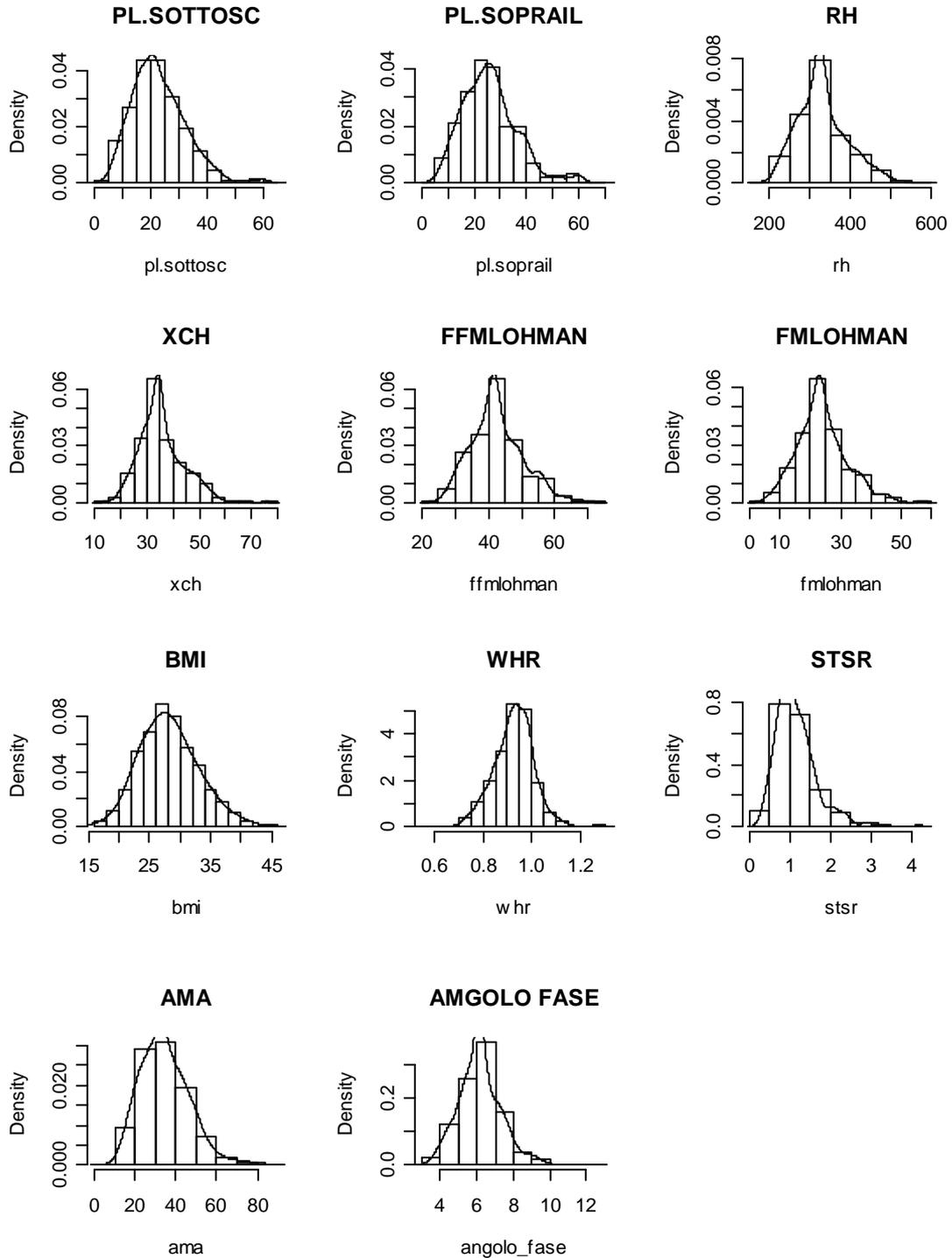
Dall'analisi dei boxplot risulta che alcune delle variabili del dataset sono asimmetriche, come `pl.bicip` e `stsr`, e questo è poco compatibile con l'ipotesi di normalità, richiesta per l'applicazione di una regressione multivariata.

Di seguito sono riportati gli istogrammi delle variabili considerate.

In R il comando per ottenere l'istogramma di una variabile è `hist(x, freq=F, main=...)`. L'opzione `freq=F` permette di ottenere, sull'asse delle ordinate, le frequenze relative anziché quelle assolute. `Main="..."` permette di indicare ad R il nome da assegnare al grafico considerato.

Con l'opzione `lines(density(x))`, R aggiunge all'istogramma delle frequenze, una stima della densità non parametrica ottenuta col metodo del nucleo (*Kernel*).

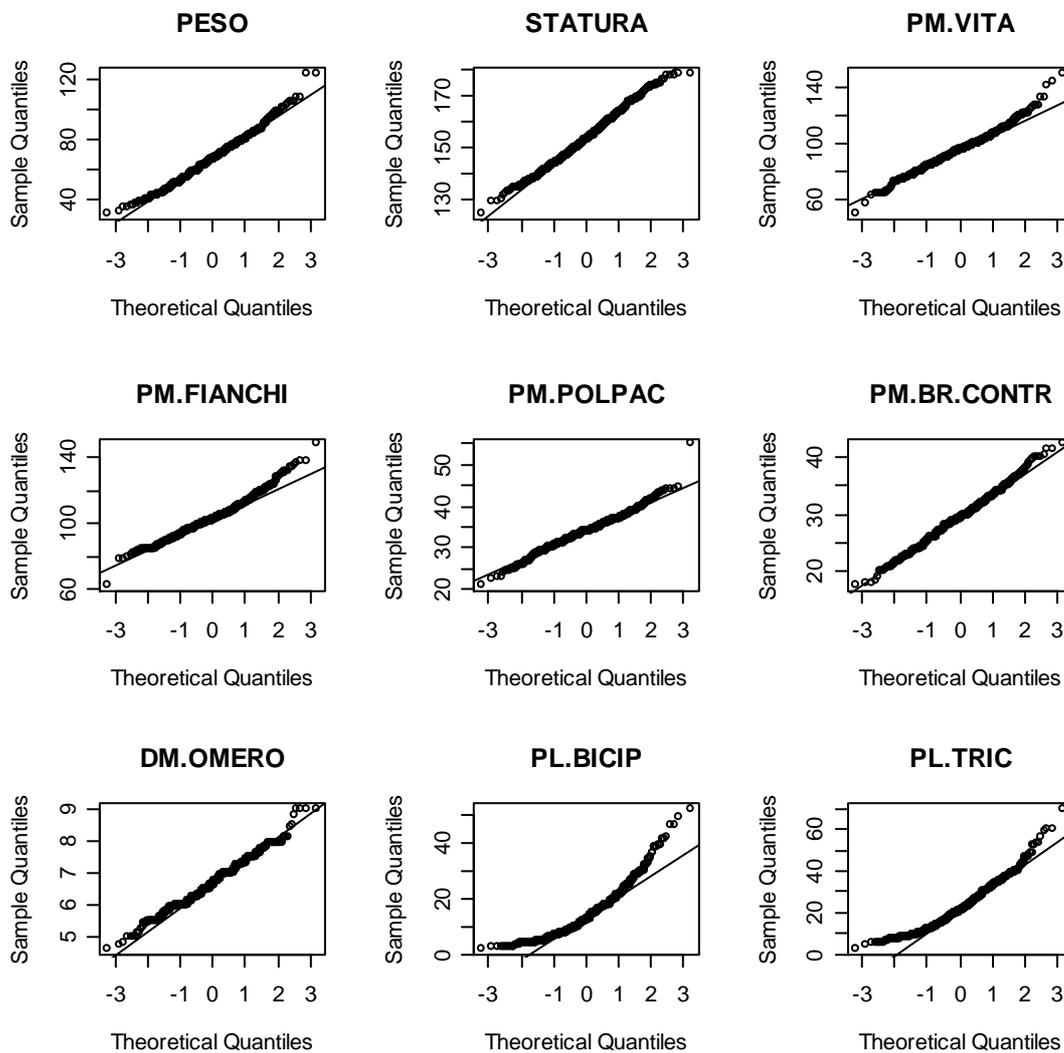


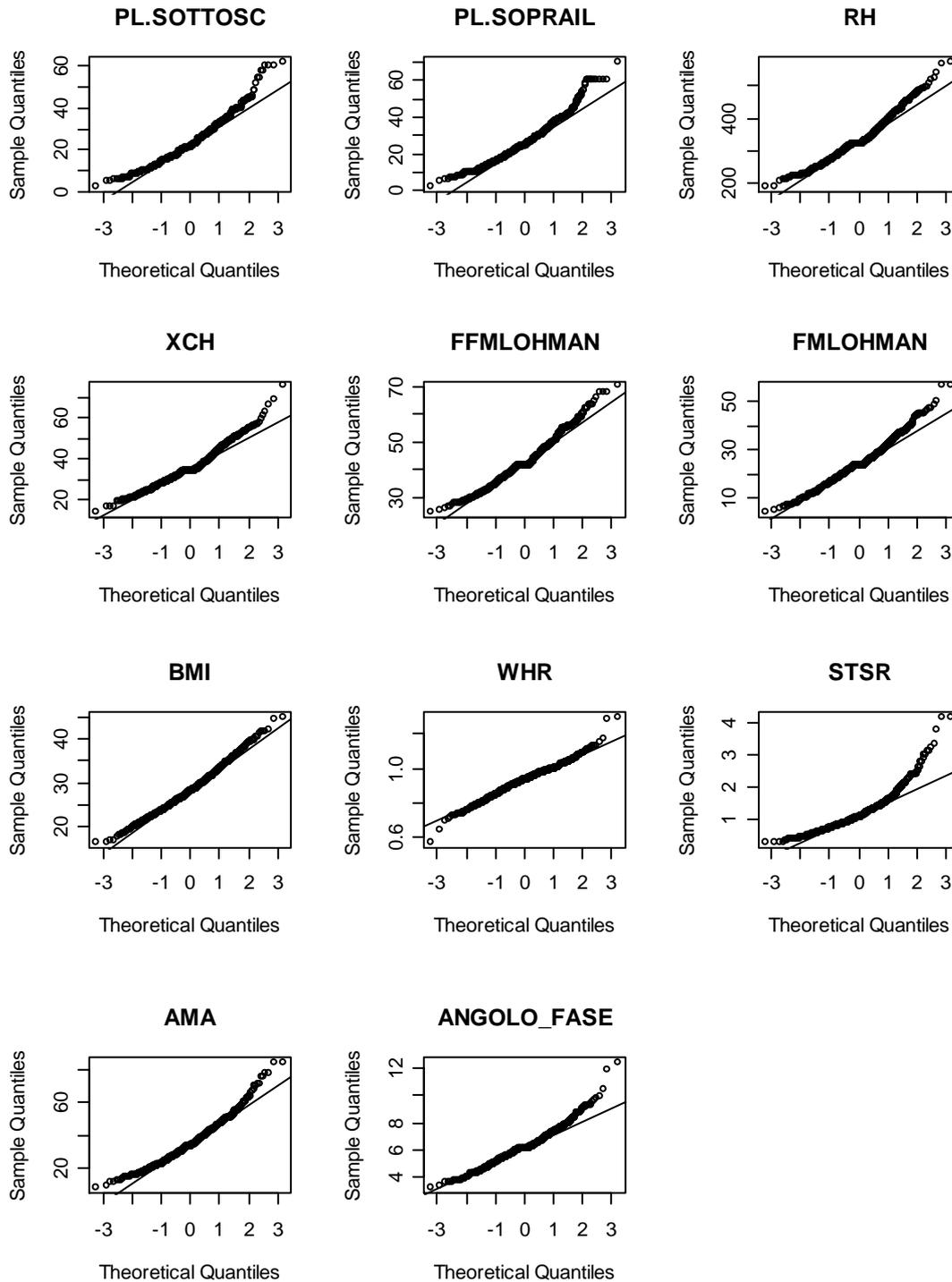


Dall'analisi degli istogrammi è possibile notare che alcune delle variabili hanno una distribuzione asimmetrica e quindi non compatibile con l'ipotesi di normalità delle

variabili risposta. In particolare le variabili `pl.bicip`, `pl.tric` e `stsr` mostrano un'asimmetria destra.

Un'ulteriore conferma del fatto che alcune delle variabili del dataset non sono normali si può ottenere utilizzando il comando `qqnorm(x)` che permette di confrontare i quantili della distribuzione della variabile di interesse con quelli teorici di una normale standard. È inoltre possibile aggiungere la retta che passa per il primo e terzo quartile della normale, per render più chiaro il confronto. Il comando per tale opzione è `qqline(x)`.



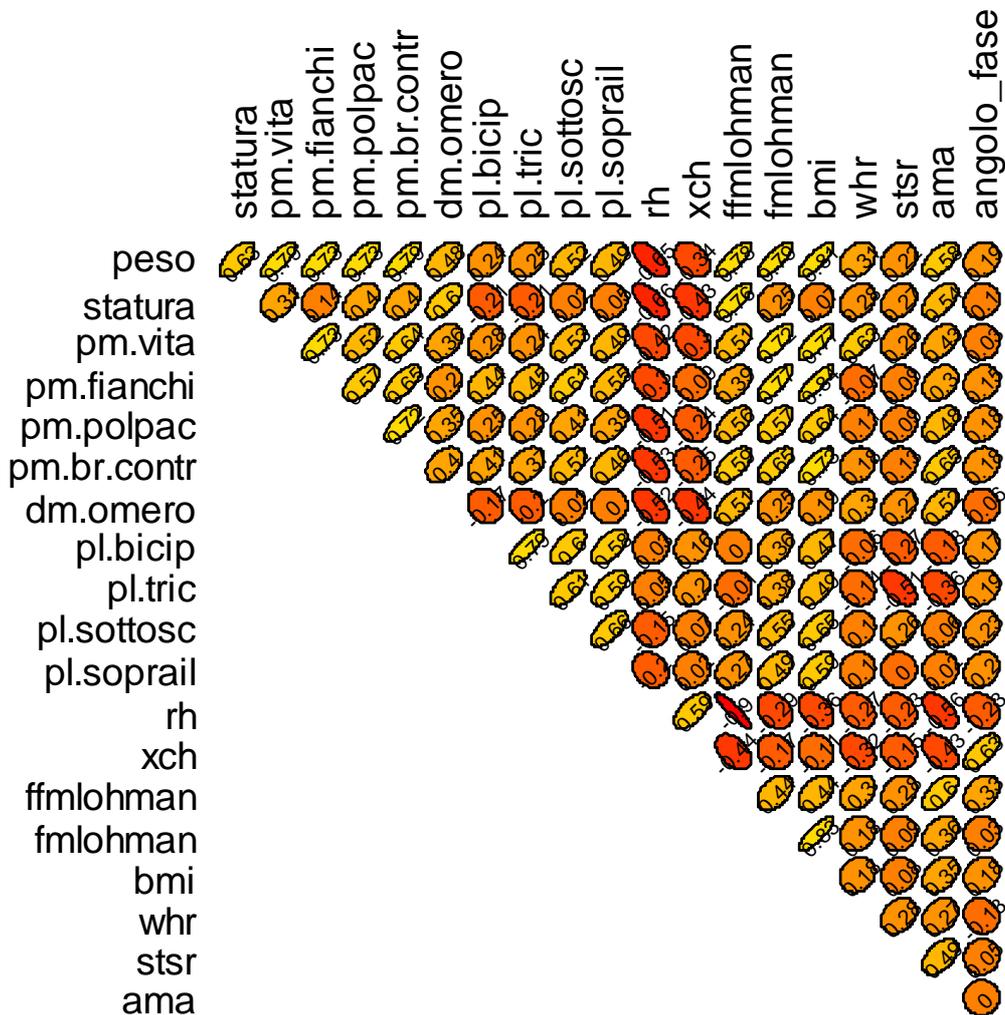


Il discostamento dei dati dalla retta fornisce la conferma definitiva che le variabili prese in esame non seguono una distribuzione normale.

In particolare risultano poco compatibili con l'ipotesi di normalità le variabili `pm.vita`, `pl.bicip`, `pl.tric`, `pl.sottosc`, `pl.sopraail`, `rh`, `stsr` ed `ama`.

3.3 Analisi bivariate

R permette, col comando `cor(data.frame)`, di calcolare le correlazioni tra tutte le variabili contenute nel dataset. Esiste inoltre la libreria “*ellipse*” che consente, tramite il comando `plotcorr(cor(...))`, di rappresentare graficamente la matrice di correlazione: dove la correlazione tra le variabili è più forte **R** rappresenta forme ellittiche; più debole è la correlazione, più la rappresentazione si avvicina a quella di una circonferenza.



Alcune delle variabili risposta sono fortemente correlate, anche perché alcune sono state calcolate come trasformazioni lineari di altre. Ad esempio il *Body Mass Index* è ottenuto

come rapporto tra peso e altezza al quadrato. Infatti la correlazione tra `bmi` e `peso` è piuttosto forte. Altre variabili del dataset sembrano incorrelate. Ad esempio l'angolo di fase mostra correlazione solo con `xch`, che rappresenta la reattanza normalizzata per la statura.

Si possono considerare anche i boxplot per le variabili risposta condizionati al sesso e alla patologia.

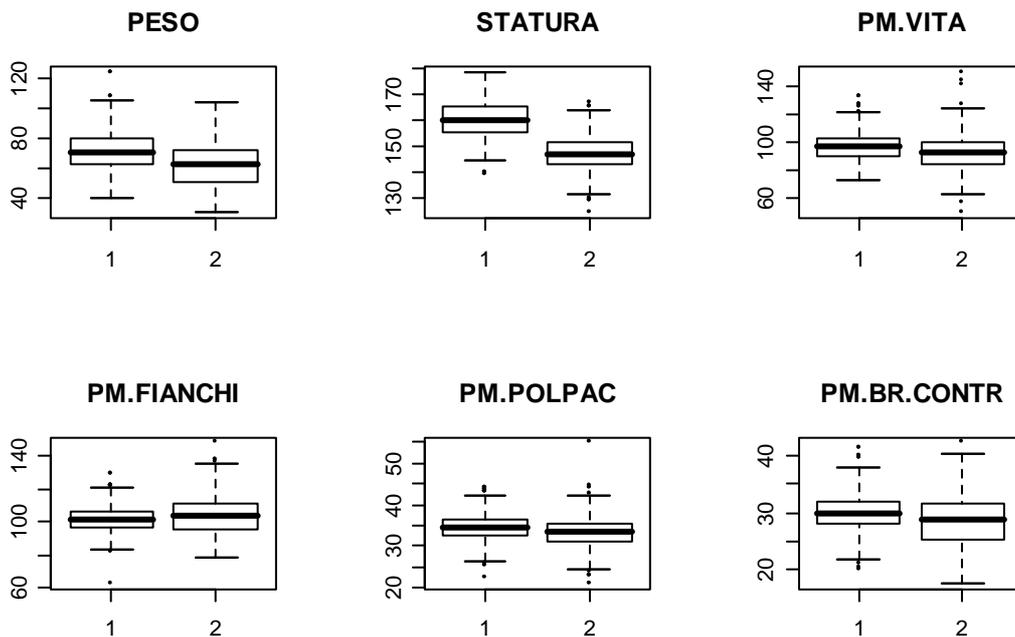
I diagrammi a scatola con baffi permettono di valutare la distribuzione di una variabile e, se condizionati ad una qualche variabile, di osservare se ci sono differenze tra le distribuzioni per i diversi valori della variabile condizionante.

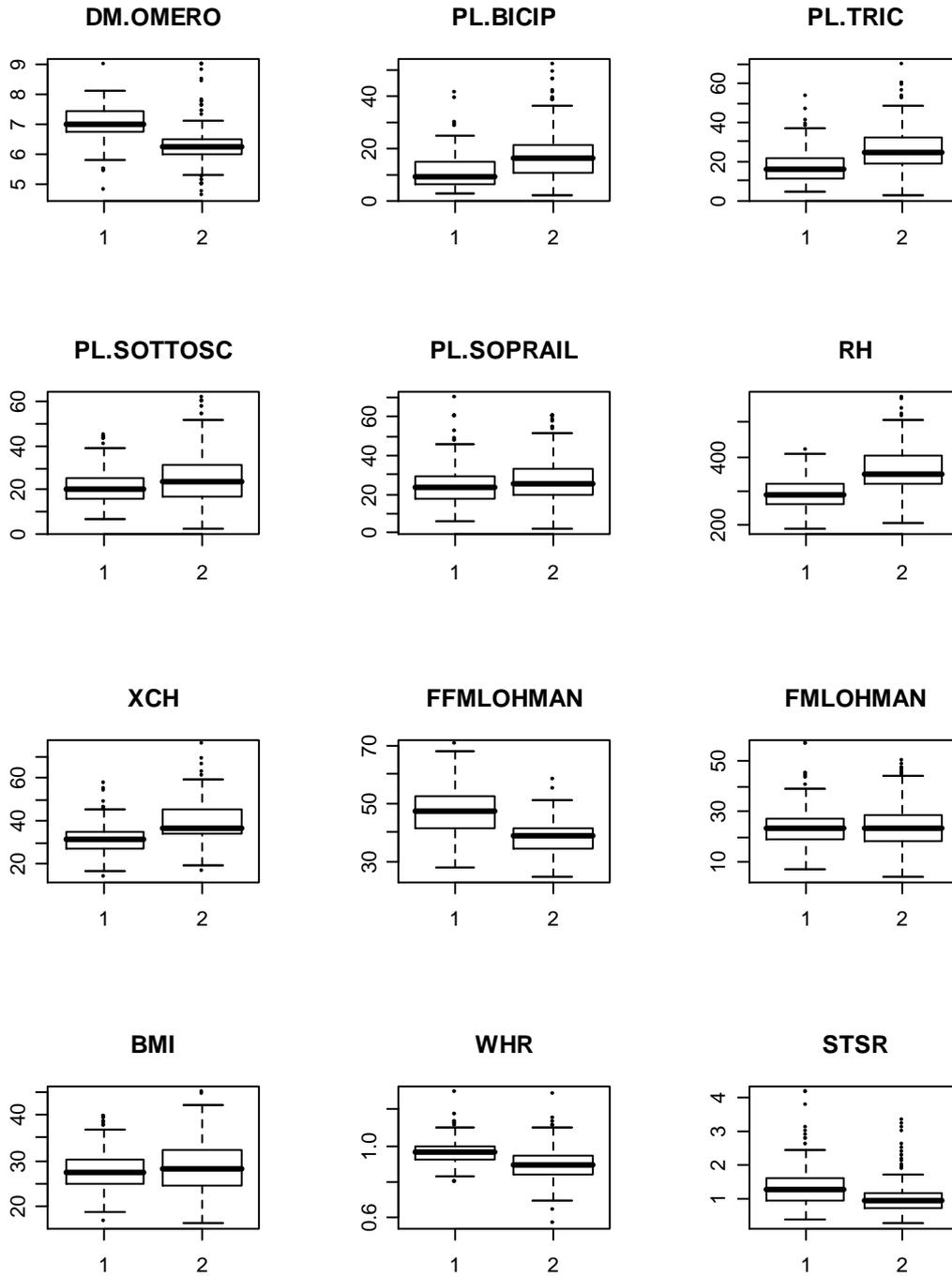
Il comando in **R** che permette di ottenere boxplot condizionati è

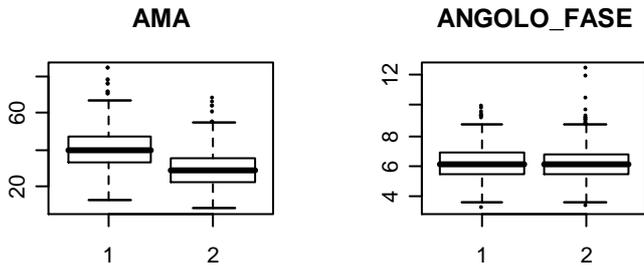
```
boxplot( y ~ x, main="..." ),
```

dove `y` è la variabile risposta e `x` è la variabile alla quale si condiziona il grafico.

Condizionatamente al sesso si ottengono i grafici riportati nel seguito.

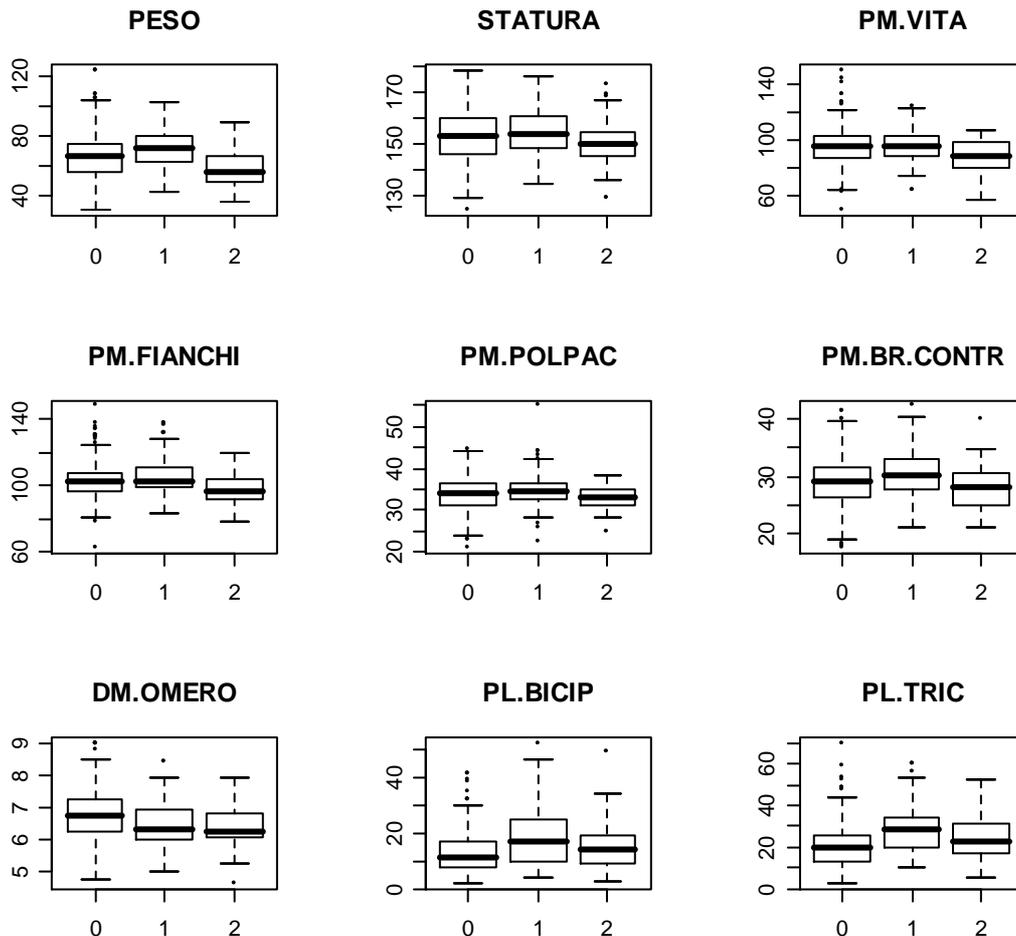


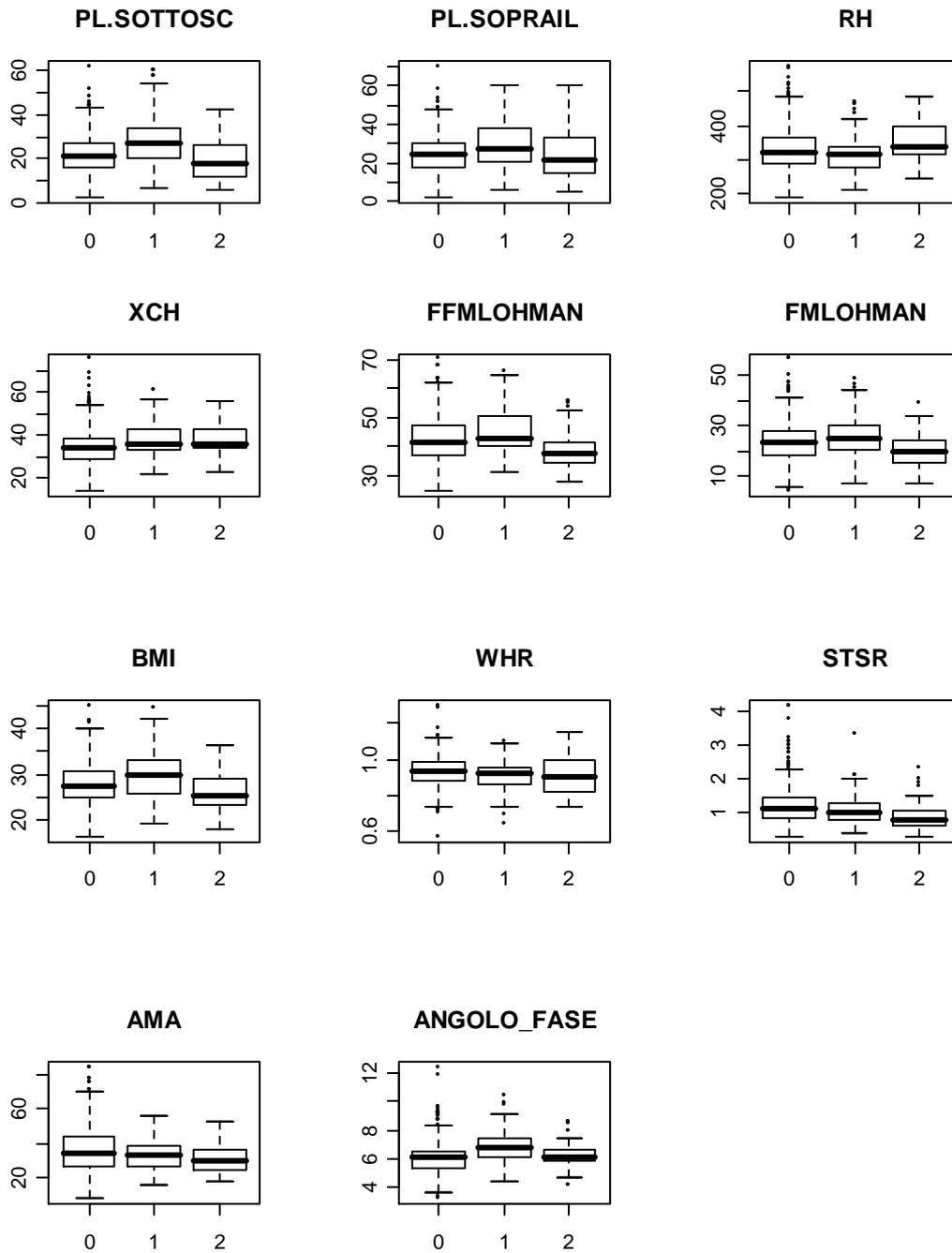




Dall'analisi dei boxplot condizionati si evince che alcune variabili assumono valori diversi a seconda del sesso, quali *dm.omero*, *stsr*, *ama* o *statura*. Altre hanno distribuzioni praticamente identiche per le due modalità, come ad esempio *angolo_fase* o *fmlohman*. Ciò significa che per alcune variabili risposta l'esplicativa sesso risulterà fortemente discriminante.

Condizionatamente alla patologia, che assume tre modalità si ottengono i seguenti grafici





Anche in questo caso si evidenzia la presenza di valori anomali nella distribuzione di alcune delle variabili che sembrano più numerosi nel caso in cui la patologia assuma

modalità 0. Le variabili `fmlohman`, `pl.sottosc` e `pl.bicip` sembrano quelle maggiormente discriminate dalla patologia.

3.4 Conclusioni

Poiché il modello di regressione multivariata si fonda sull'ipotesi di normalità delle risposte si possono provare delle trasformazioni delle variabili che le rendano, approssimativamente, normali. Trasformazioni possibili sono le funzioni logaritmo e radice quadrata, o la trasformata di Box-Cox .

La trasformata di Box-Cox assume la seguente forma

$$\begin{cases} \frac{Y^\lambda - 1}{\lambda} \\ \log(Y) \end{cases}$$

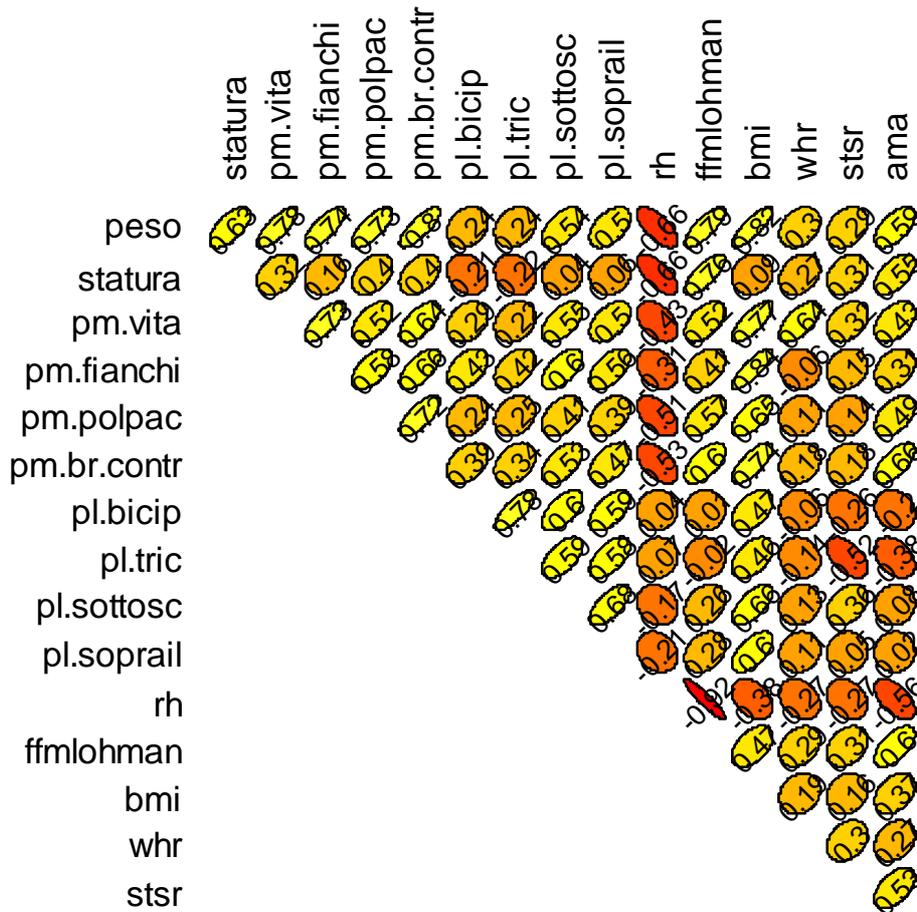
ed in particolare è pari a $\log(Y)$ quando $\lambda = 0$. Inoltre, anche la radice quadrata si può pensare come un caso particolare della trasformata di Box-Cox quando $\lambda = 1/2$. Tale trasformata è particolarmente utile a stabilizzare la varianza (si veda Di Fonzo-Lisi, 2005 [3]).

Dalle suddette analisi si conclude che il dataset definitivo, che verrà usato per stimare il modello è costituito dalle seguenti variabili:

- `sqrt(peso)` ;
- `log(statura)` ;
- `sqrt(pm.vita)` ;
- `log(pm.fianchi)` ;
- `pm.polpac` ;
- `pm.br.contr` ;
- `BoxCox(pl.bicip, 0.4)` ;
- `BoxCox(pl.tric, 0.4)` ;
- `BoxCox(pl.sottosc, 0.4)` ;
- `sqrt(pl.soprail)` ;
- `BoxCox(rh, 0.7)` ;
- `sqrt(ffmlohman)` ;
- `log(bmi)` ;

- whr ;
- log(stsr) ;
- sqrt(ama) .

Le variabili `dm.omero`, `xch`, `fmlohman` e `angolo_fase` non si possono assumere normali, neanche se trasformate. Per questo verranno escluse nelle analisi di seguito riportate. Può essere utile valutare le correlazioni tra le nuove variabili create



Sul nuovo dataset sarà stimato un modello di regressione multivariata, avente come risposte le variabili trasformate. Prima di stimare il modello sui dati è possibile cercare di valutare l'effetto marginale delle variabili esplicative sulle risposte. Per quanto riguarda il sesso e la patologia, sono già state effettuate alcune analisi in proposito tramite i boxplot condizionati. È, quindi possibile valutare l'effetto del sesso sulle variabili riposta tramite un test t a due campioni; l'ipotesi nulla del test è che le medie

delle due popolazioni, quella maschile e quella femminile, per la variabile considerata, siano uguali. Il comando usato in **R** per ottenere un test *t* a due campioni è

`t.test(x, ...)`, dove *x* è il vettore dei dati.

Di seguito è riportata la tabella contenente il valore del *t*-test e il relativo *p-value*.

Tutte le analisi successive sono state effettuate sulle variabili trasformate in modo da essere compatibili con la normalità.

***t*-test (sesso)**

<u>Variabile</u>	<u>statistica test</u>	<u>p-value</u>
<u>peso</u>	10.3521	< 2.2e-16
<u>statura</u>	26.9829	< 2.2e-16
<u>pm.vita</u>	5.719	1.561e-08
<u>pm.fianchi</u>	-2.7481	0.006147
<u>pm.polpac</u>	4.683	3.363e-06
<u>pm.br.contr</u>	4.7543	2.400e-06
<u>pl.bicip</u>	-11.8423	< 2.2e-16
<u>pl.tric</u>	-13.6355	< 2.2e-16
<u>pl.sottosc</u>	-4.3024	1.924e-05
<u>pl.soprail</u>	-2.7935	0.005349
<u>rh</u>	-20.6872	< 2.2e-16
<u>ffmlohman</u>	19.5933	< 2.2e-16
<u>bmi</u>	-1.3679	0.1718
<u>whr</u>	12.6781	< 2.2e-16
<u>stsr</u>	11.1546	< 2.2e-16
<u>ama</u>	13.9487	< 2.2e-16

Il sesso non risulta discriminante solo per la variabile *bmi*. Per tutte le altre variabili risposta si rifiuta l'ipotesi che ci sia uguaglianza in media nei due gruppi.

Per quanto riguarda la variabile *patologia*, è possibile effettuare il `oneway.test` che permette di verificare l'ipotesi di uguaglianza in media di più di due gruppi.

Tale test consiste in un'analisi della varianza ed è una generalizzazione del *t*-test.

Di seguito si riportano i valori della statistica test e del *p-value*.

<u>Variabile</u>	<u>Valore del test F</u>	<u>p-value</u>
<u>peso</u>	22.224	4.42e-09
<u>statura</u>	5.4096	0.005493
<u>pm.vita</u>	8.5718	0.0003142
<u>pm.fianchi</u>	11.6249	2.248e-05
<u>pm.polpac</u>	7.3813	0.0008974
<u>pm.br.contr</u>	9.1997	0.0001799
<u>pl.bicip</u>	24.5241	1.069e-09
<u>pl.tric</u>	45.2213	1.059e-15
<u>pl.sottosc</u>	27.829	9.655e-11
<u>pl.soprail</u>	10.835	4.676e-05
<u>rh</u>	9.2292	0.0001753
<u>ffmlohman</u>	13.588	4.231e-06
<u>bmi</u>	16.9252	2.825e-07
<u>whr</u>	4.4189	0.01394
<u>stsr</u>	15.7685	7.42e-07
<u>ama</u>	6.6231	0.001752

Al 5% di significatività si rifiuta l'ipotesi nulla di uguaglianza della media nei tre gruppi per tutte le variabili. All'1% si accetterebbe solo per la variabile *whr*.

Per valutare l'effetto marginale dell'età, che è una variabile quantitativa, si può analizzare la correlazione di *eta* con le variabili risposta.

	<u>peso</u>	<u>statura</u>	<u>pm.vita</u>	<u>pm.fianchi</u>	<u>pm.polpac</u>	<u>pm.br.contr</u>
<u>cor</u>	-0.471	-0.405	-0.169	-0.222	-0.458	-0.498

	<u>pl.bicip</u>	<u>pl.tric</u>	<u>pl.sottosc</u>	<u>pl.soprail</u>	<u>rh</u>	<u>ffmlohman</u>
<u>cor</u>	-0.078	-0.115	-0.298	-0.263	0.371	-0.473

	<u>bmi</u>	<u>whr</u>	<u>stsr</u>	<u>ama</u>
<u>cor</u>	-0.313	0.011	-0.178	-0.389

Si può notare che la correlazione assume valore negativo per la maggior parte delle variabili, eccetto che per r_h e w_{hr} . Inoltre la correlazione tra eta e w_{hr} assume un valore piuttosto piccolo. La maggior parte delle correlazioni assume un valore maggiore di 0.2 in valore assoluto. Questo potrebbe essere un segnale significativo della presenza di un effetto marginale della variabile relativa all'età sulle risposte. In seguito a queste analisi si può effettuare la stima del modello, tenendo anche conto degli effetti congiunti delle variabili, al fine di valutarne la significatività.

Capitolo 4

Adattamento ai dati

4.1 Il pacchetto *lm*

Il comando *lm* in R è utilizzato per stimare modelli lineari, sui quali effettuare regressioni, analisi della varianza e della covarianza (in particolare l'opzione `aov` viene utilizzata per questi scopi).

L'espressione del comando è:

```
lm(formula, data, subset, ...)
```

Argomenti:

- `formula` indica la specificazione del modello che si vuole stimare;
- `data` indica il dataset sul quale effettuare la regressione espressa nella formula e contiene le variabili del modello; se l'opzione viene omessa, **R** effettua la regressione sul dataset sul quale è stato effettuato l'*attach*;
- `subset` si utilizza per selezionare un sottoinsieme di osservazioni sul quale effettuare la regressione.

La formula del modello, nel comando *lm*, è del tipo “risposte ~ esplicative”, dove “risposte” indica il vettore o la matrice delle variabili dipendenti del modello, mentre “esplicative” indica la matrice delle variabili indipendenti, che, combinate linearmente, modellano la risposta. Se, ad esempio, nella formula abbiamo due esplicative, x_1 e x_2 , e

il modello è espresso come $y \sim x_1 + x_2$, ciò vuol dire che sulla variabile risposta si considera l'effetto di x_1 e l'effetto di x_2 separatamente.

Infine se il modello è del tipo $y \sim x_1 * x_2$, si valutano, sulla risposta y , gli effetti di x_1 e x_2 presi separatamente e l'effetto della loro interazione; $x_1 * x_2$ corrisponde ad $x_1 + x_2 + x_1:x_2$.

Nel caso in cui la formula sia del tipo $y \sim x - 1$, ciò significa che si sta considerando il modello di regressione senza intercetta.

Un caso particolare della regressione col comando `lm` si ha quando Y è una matrice ($n \times p$) di risposte e non un vettore; questo è il caso di interesse nel caso di regressione multipla multivariata.

In questa situazione R stima separatamente le regressioni ai minimi quadrati per ogni colonna della matrice Y .

Utilizzando il comando `summary()` si ottiene un sommario dell'output, comprendente le stime e gli *standard error* dei coefficienti, con i corrispondenti valori dei test per la verifica della significatività e i relativi *p-value*. Inoltre vengono visualizzati il minimo, il massimo, il primo e il terzo quartile e la mediana della distribuzione dei residui del modello.

Infine compaiono anche i valori dell' R^2 e dell' R^2 corretto, il *residual standard error*, che rappresenta la quantità di varianza che non è spiegata dal modello, con i relativi gradi di libertà e la statistica F per il modello con il solo termine d'intercetta, con il *p-value* ad essa associato.

Una volta stimato il modello di regressione multivariata, si può utilizzare la funzione `anova()` per confrontare modelli "annidati" e per valutare se sia opportuno togliere o aggiungere una determinata variabile esplicativa nel modello; ciò si effettua considerando se risulta significativo il peggioramento o il miglioramento in termini di varianza residua, che si ottiene con la modifica adottata.

4.1.1 Il comando `anova.mlm()`

Nel caso di modelli multivariati, l'estensione dell'opzione `anova()` è il comando `anova.mlm()`.

L'espressione di tale comando è:

```
anova.mlm(object,..., test = c("Pillai", "Wilks", "Hotelling-Lawley",
"Roy", "Spherical"),...)
```

dove:

- `object` è un oggetto della classe `mlm`;
- `test` indica il tipo di test che si desidera visualizzare in output: Λ di Wilks, T^2 di Hotelling, Pillai's *Trace*, etc..

4.2 La stima del modello completo

Il modello più complesso che si può stimare sui dati, presentati nel capitolo precedente, inserisce tra le esplicative sesso, patologia, età e tutte le interazioni possibili tra le tre variabili.

Il comando per stimare con R questo modello è il seguente

```
fit<-lm(cbind(peso, statura, pm.vita, pm.fianchi, pm.polpac, pm.br.contr,
pl.bicip, pl.tric, pl.sottosc, pl.soprail, rh, ffmlohman, bmi, whr, stsr,
ama) ~ sesso * patologia * eta)
```

Con la stima di questo modello si ottiene una matrice \hat{B} di coefficienti stimati di dimensioni (11×16) , perché le variabili risposta sono 16 e le esplicative, poiché si considerano tutte le interazioni possibili, sono 11.

Più che valutare il valore dei singoli coefficienti, ha senso cercare di determinarne la significatività, per stabilire se sia possibile semplificare il modello.

Si può inoltre valutare, in termini di devianza residua, la bontà del modello e il grado di adattamento ai dati.

Inoltre per tutte le variabili si rifiuta l'ipotesi che sia buono il modello con la sola intercetta come variabile esplicativa.

L'intercetta è l'unico parametro del modello significativo, a tutti i livelli di significatività usualmente considerati, per tutte le variabili.

L'età ha un effetto rilevante e negativo per la maggior parte delle risposte, eccetto che per `pl.bicip`, `pl.tric` ed `whr` per le quali il parametro risulta non significativo e `rh`, che è l'unica variabile per la quale la stima del coefficiente di `eta` è positivo.

Poiché gli antropologi sono interessati alla costruzione di curve di crescita, è importante notare che l'età è una variabile esplicativa importante per spiegare gli indicatori antropometrici calcolati e che ha una relazione negativa con questi ultimi: l'aumentare dell'età porta, indicativamente, ad una riduzione delle quantità considerate.

La variabile `sezzo` sembra avere un effetto particolarmente significativo sulle variabili `statura`, `pm.vita`, `pl.bicip`, `pl.tric`, `pl.sottosc`, `rh`, `ffmlohman`, `whr`, ed `ama`. Il *p-value* associato all'ipotesi di nullità del coefficiente della variabile `sezzo`, per queste variabili, è inferiore a 0.05.

I coefficienti associati alle due variabili dicotomiche che codificano la patologia sono significativi, entrambi o almeno uno dei due, al 5% per le seguenti variabili: `pl.bicip`, `pl.tric`, `pl.sottosc` ed `ama`.

Riguardo alle interazioni tra le variabili esplicative del modello, l'interazione a tre è significativa solo per le variabili `pl.tric`, `rh` e `ffmlohman`. L'interazione tra età e sesso è rilevante per spiegare `pm.vita`, `pl.bicip`, `pl.tric`, `pl.sottosc`, `ffmlohman` e `whr`.

L'interazione tra patologia e sesso ha un coefficiente significativo al 5% per `pl.tric`, `rh`, `ffmlohman` e `stsr`.

Infine l'interazione tra patologia ed età è indicativa per le variabili `pl.bicip`, `pl.tric`, `pl.sottosc` ed `ama`.

```
anova.mlm(fit)
```

Analysis of Variance Table

	Df	Pillai approx	F num	Df den	Df	Pr(>F)
(Intercept)	1	1 24071478		16	716	< 2.2e-16 ***
sezzo	1	1	116	16	716	< 2.2e-16 ***
patologia	2	0.30496	8	32	1434	< 2.2e-16 ***
eta	1	0.42672	33	16	716	< 2.2e-16 ***
sezzo:patologia	2	0.13195	3	32	1434	9.412e-09 ***
sezzo:eta	1	0.20541	12	16	716	< 2.2e-16 ***
patologia:eta	2	0.09351	2	32	1434	0.0001422 ***

```

sesso:patologia:eta    2  0.06736          2    32   1434 0.0241187 *
Residuals              731

```

Dall'analisi dell'output del comando `anova.mlm()` sul modello `fit`, risulta che eliminando il termine di interazione a tre dalla regressione multipla, il peggioramento in termini di devianza residua che si otterrebbe non è significativo all'1% ($p\text{-value} = 0.02$). Si può provare ad eliminare l'interazione a tre e vedere come cambia il modello.

4.2.1 Semplificazione del modello

Per eliminare il termine che rappresenta l'interazione a tre tra le esplicative del modello si usa il comando

```
fit1<-update(fit, ~.-sesso:patologia:eta)
```

dove il comando `update` indica che si sta aggiornando, modificando il modello `fit`, escludendo il termine `sesso:patologia:eta`.

La nuova matrice di coefficienti stimati è ora di dimensioni 9×16 .

I *residual standard errors* del nuovo modello sono poco differenti da quelli del modello precedente. Inoltre in questo caso vanno confrontati con 733 e non 731 gradi di libertà. Se valutato per ogni singola variabile risposta, il modello più semplice sembra essere buono tanto quanto il precedente.

```

anova(fit1)
Analysis of Variance Table

              Df    Pillai approx F num Df den Df    Pr(>F)
(Intercept)    1         1 23875740    16   718 < 2.2e-16 ***
sesso           1         1     116     16   718 < 2.2e-16 ***
patologia       2    0.30144         8     32  1438 < 2.2e-16 ***
eta             1    0.42497         33     16   718 < 2.2e-16 ***
sesso:patologia 2    0.13153         3     32  1438 9.586e-09 ***
sesso:eta       1    0.20529         12     16   718 < 2.2e-16 ***
patologia:eta   2    0.09311         2     32  1438 0.0001470 ***

```

Residuals 733

I termini inseriti nel modello sembrano tutti ampiamente significativi. Se si confrontano i due modelli si ottiene

```
anova.mlm(fit,fit1)
Analysis of Variance Table
  Res.Df  Df Gen.var.  Pillai approx F num Df den Df  Pr(>F)
1     731      0.14217
2     733     2  0.14240  0.06736  1.56199     32  1434  0.02412 *
```

Da ciò si evidenzia che mentre per le singole variabili il termine di interazione a tre tra sesso, patologia ed eta non era significativo, la sua eliminazione dal modello comporta una perdita significativa nel modello complessivo, nel quale si tiene conto anche della correlazione tra le variabili *Y*.

Se si eliminasse anche una interazione a due, il peggioramento del modello risulterebbe ancora significativo. Si possono stimare tre modelli

```
fit21<-update(fit1,~.-patologia:eta)
fit22<-update(fit1,~.-sesso:eta)
fit23<-update(fit1,~.-sesso:patologia)
```

e confrontandoli con *fit1* si ottiene

```
anova.mlm(fit1,fit21)
Analysis of Variance Table
  Res.Df  Df Gen.var.  Pillai approx F num Df den Df  Pr(>F)
1     733      0.14240
2     735     2  0.14286  0.09311  2.19420     32  1438  0.0001470 ***
```

```
anova.mlm(fit1,fit22)
Analysis of Variance Table
  Res.Df  Df Gen.var.  Pillai approx F num Df den Df  Pr(>F)
1     733      0.1424
2     734     1  0.1443  0.2095  11.8914     16  718  < 2.2e-16 ***
```

```
anova.mlm(fit1,fit23)
```

```
Analysis of Variance Table
```

	Res.Df	Df	Gen.var.	Pillai	approx F	num Df	den Df	Pr(>F)
1	733		0.14240					
2	735	2	0.14304	0.11119	2.64535	32	1438	2.007e-06 ***

Si conclude quindi che il modello più significativo per spiegare i dati è quello completo.

Si può quindi, stimare sui dati il modello ritenuto opportuno.

Dall'analisi del modello ridotto risulta che il coefficiente dell'intercetta è significativo per tutte e sedici le variabili risposta.

Il coefficiente della variabile dicotomica *sex* risulta significativo al 5% per le variabili *pm.vita*, *pl.tric*, *pl.sottosc*, *rh*, *ffmlohman*, *whr*, *statura* e *ama*. Per le ultime due, però, il coefficiente risulta non significativo se si fissa come livello d'accettazione l'1%.

La variabile *patologia* risulta significativa al 5% nelle regressioni delle seguenti variabili: *peso*, *statura*, *pm.fianchi*, *pm.polpac*, *pl.sottosc*, *bmi*, *ama*, *pm.vita* e *pm.br.contr*. Per le ultime due variabili, la *patologia* non risulta avere un effetto significativo se il livello considerato è l'1%. Inoltre al 10% la *patologia* è significativa anche per la variabile *ffmlohman*.

Per quanto riguarda la variabile *età*, il suo effetto è significativo su tutte le risposte al 5% di significatività, eccetto che su *pl.bicip*, *pl.tric* e *whr*. Inoltre, se il livello di accettazione è l'1%, l'età non risulta rilevante neanche nella regressione della variabile *stsr*.

Oltre all'effetto delle singole esplicative, si può valutare se le interazioni tra esse risultano significative.

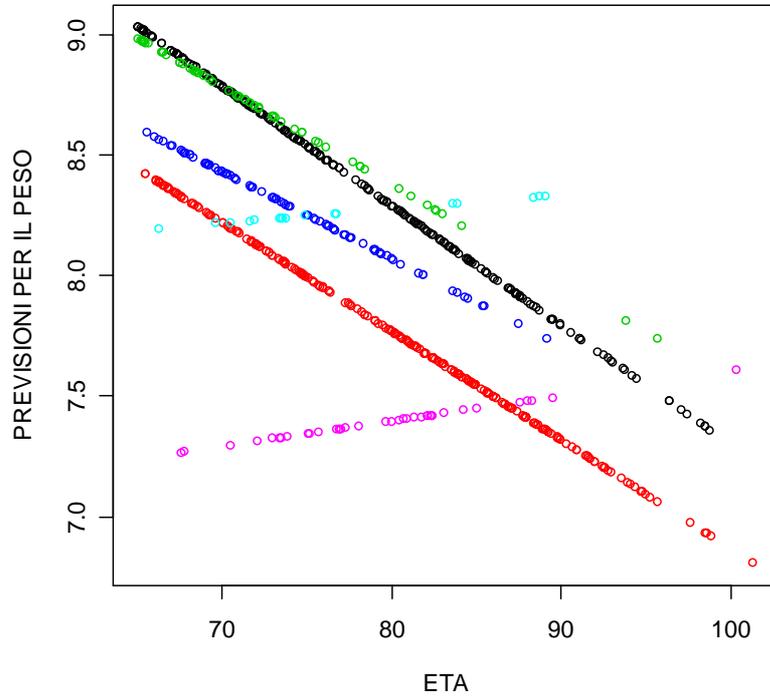
L'interazione tra *sex* ed *età*, ha un effetto significativo al 5% solo sulle variabili *pm.vita*, *pl.bicip*, *pl.tric*, *pl.sottosc*, *ffmlohman* e *whr*; i coefficienti, per altro, risultano significativi all'1% solo per le ultime due variabili. Tenendo conto, però, delle correlazioni tra le risposte, complessivamente il coefficiente dell'interazione tra *sex* ed *età* risulta rilevante.

I coefficienti dell'interazione tra sesso e patologia sono significativi, entrambi o uno dei due, al 5% per le variabili `pm.vita`, `pm.polpac`, `pm.br.contr`, `pl.sopraail`, `stsr`, `ama`, `pm.fianchi`, `pl.sottosc` e `whr`. In particolare, per le ultime tre variabili il coefficiente dell'interazione non risulterebbe significativo se il livello di significatività considerato fosse l'1%. Al 10% l'interazione tra sesso e patologia è rilevante anche nella regressione per `bmi`.

Infine, l'interazione tra patologia e età è significativa al 5% nelle regressioni delle seguenti risposte: `ama`, `peso`, `statura`, `pm.fianchi`, `pm.polpac`, `pm.br.contr`, `pl.sottosc` e `bmi`. Nelle regressioni relative alle ultime tre variabili, il coefficiente di interazione non risulterebbe significativo all'1%. Se fa riferimento ad una probabilità del 10%, risultano significativi anche i coefficienti dell'interazione nelle regressioni delle variabili `pm.vita`, `rh` e `ffmlohman`.

4.3 Previsioni

A partire dal modello considerato è possibile costruire i valori predetti con l'opzione `fitted(fit)` del comando `lm`. Con tale comando si ottiene una matrice di previsioni di dimensioni (743×16) poiché nel caso in esame la regressione effettuata è multivariata. Si può, quindi, dalla matrice estrarre il vettore delle previsioni per una variabile di interesse. In particolare, le previsioni costruite possono essere rappresentate graficamente in un diagramma di dispersione in cui la variabile esplicativa è l'età. Inoltre, è possibile distinguere in tali grafici le unità statistiche in base al sesso e alla patologia rilevata. Discriminando in base ai due fattori, le curve di crescita che si ottengono per le previsioni della variabile risposta `peso` sono le seguenti



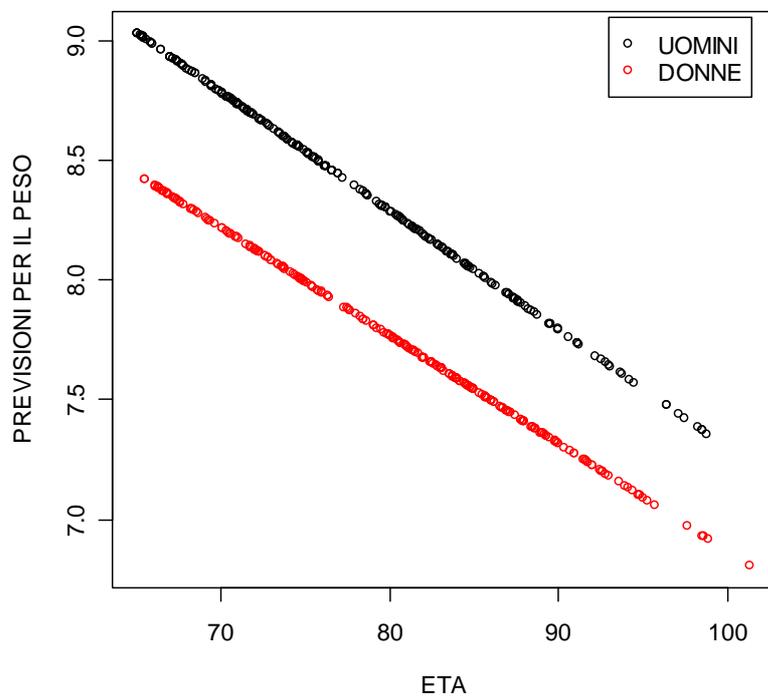
LEGENDA	
◊ uomini sani	◊ donne diabetiche
◊ donne sane	◊ uomini con Alzheimer
◊ uomini diabetici	◊ donne con Alzheimer

Si ricorda che le previsioni per la variabile peso sono ottenute a partire dalla radice quadrata della variabile contenuta nel dataset di partenza, effettuata ai fini della normalizzazione. A causa di ciò i valori in ordinata sul grafico, sono compresi tra 7.0 e 9.0, circa.

Come si può osservare, per ognuna delle due modalità del sesso si distinguono tre differenti curve, che corrispondono alla tre modalità della variabile patologia.

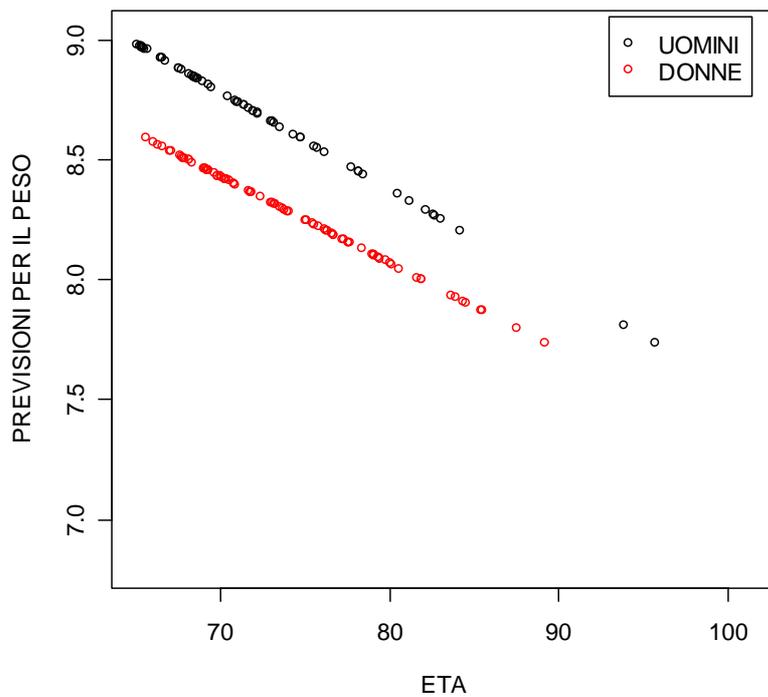
Se, invece, si rappresentano, separatamente per la patologia, le curve relative ai maschi e alle femmine, si ottiene, per gli individui sani, il seguente grafico

INDIVIDUI SANI

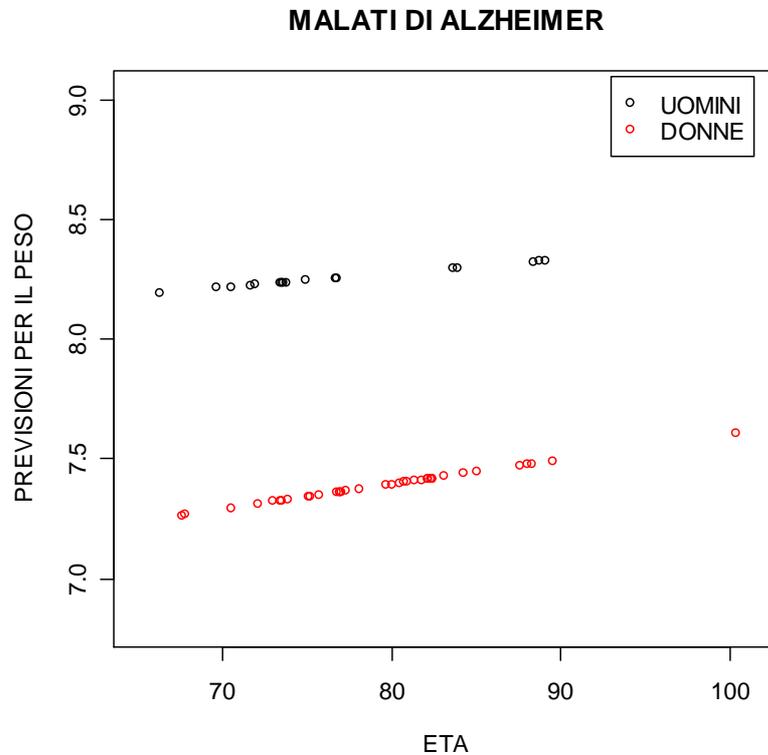


Per i diabetici il grafico ottenuto è il seguente

DIABETICI



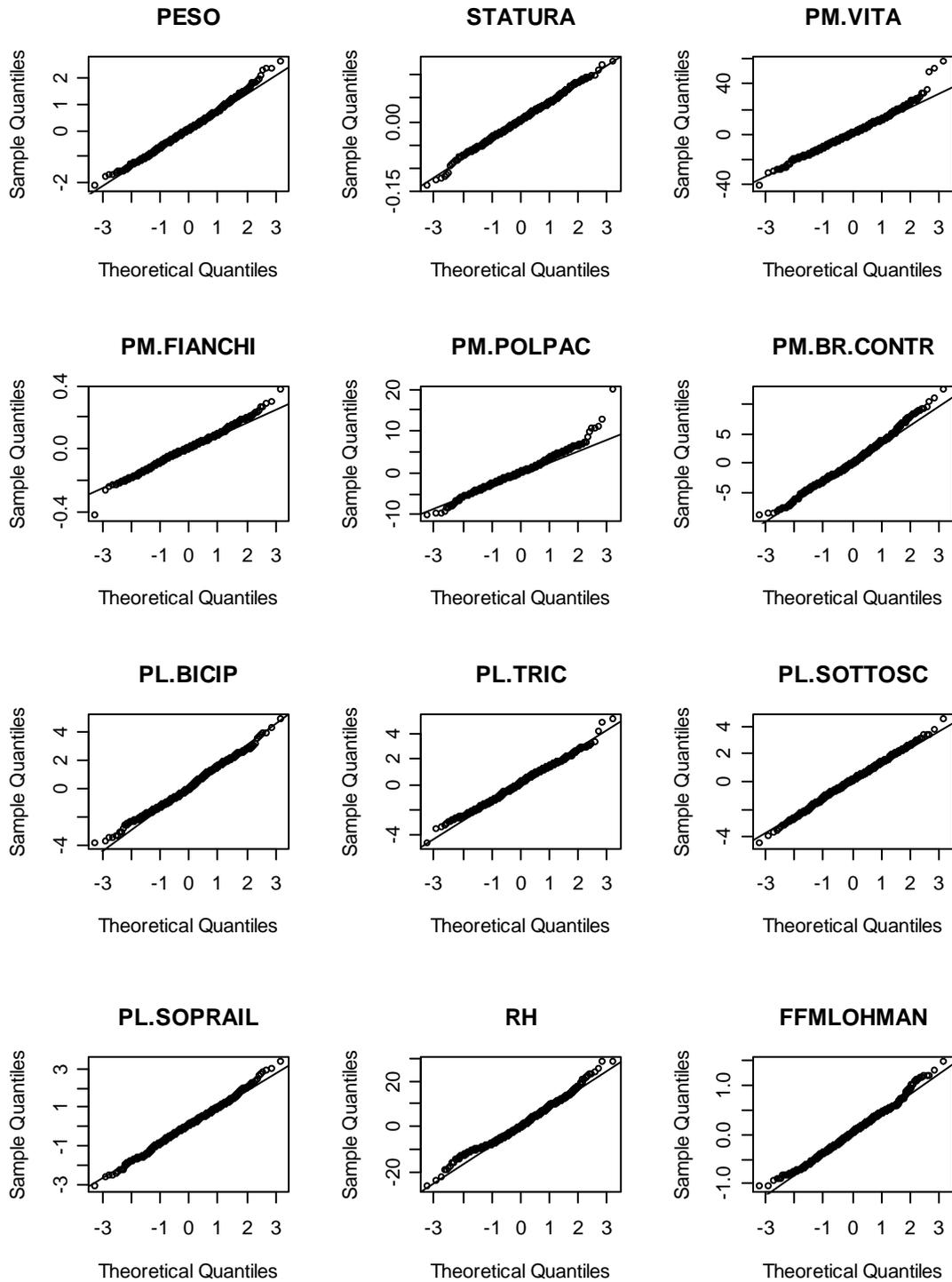
Infine, per i malati di Alzheimer si ha

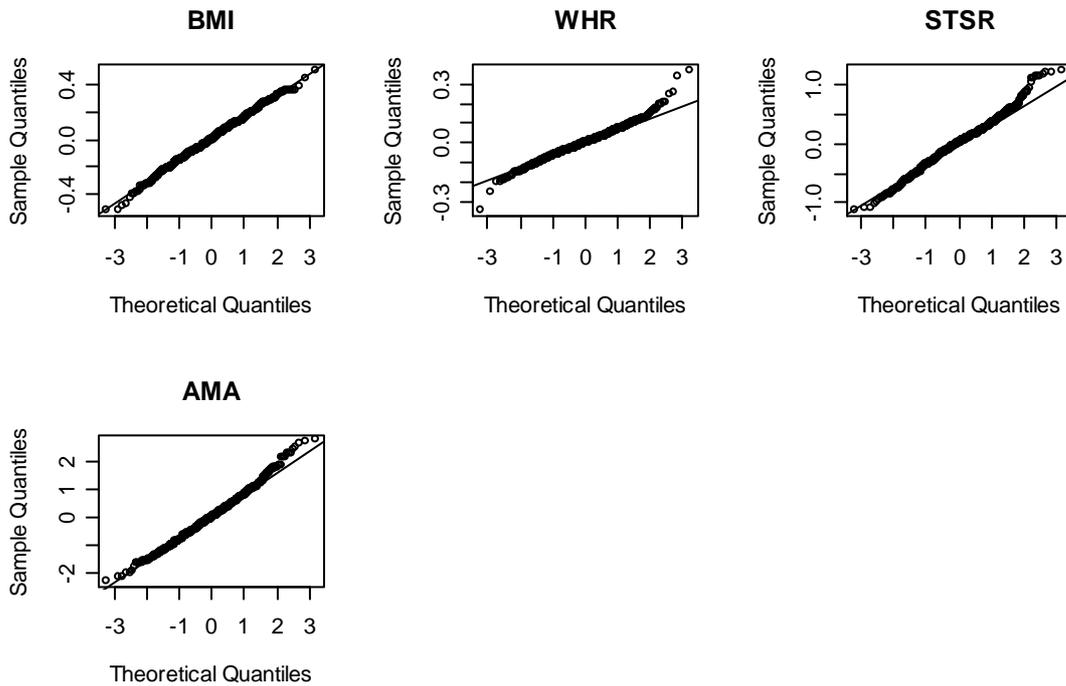


Come si può notare, le curve sono decrescenti per gli individui sani e per i diabetici, appaiono invece costanti o al più debolmente crescenti per i malati di Alzheimer. Il peso, dunque, sembra diminuire con l'età per i primi due gruppi di individui, a prescindere dal sesso. Per i malati di Alzheimer, invece, il peso sembra rimanere praticamente costante al variare dell'età. Inoltre, mentre per gli individui sani e per i malati di Alzheimer le curve distinte per i due sessi sembrano parallele, per i diabetici sembra esserci un cambiamento di pendenza.

4.4 Conclusioni

Il modello stimato può essere validato anche a partire dall'analisi dei residui, che si possono calcolare come differenza tra valori osservati e valori predetti. Secondo le ipotesi alla base del modello, i residui dovrebbero provenire da una distribuzione normale. È quindi possibile valutare graficamente se provengono da una normale, attraverso i *q-qplot*.





I grafici quantile-quantile sembrano compatibili con l'ipotesi di normalità per tutte le variabili, eccetto che per *pm.vita*, *pm.polpac*, *whr* e *stsr* per le quali potrebbero esserci dei dubbi.

Dalle analisi effettuate in seguito alla stima del modello senza interazione a tre, si può concludere che, benché valutati per ogni singola variabile risposta, talvolta gli effetti delle variabili esplicative del modello sembrano non essere rilevanti, tuttavia, tenendo conto della correlazione delle risposte, che assume per alcune valori significativamente alti, tutti i coefficienti stimati nel modello sono significativi.

Ciò porta a concludere che le variabili esplicative considerate siano degli indicatori adeguati a spiegare le variabili risposta considerate. In particolare, le variabili di natura antropometrica sembrano avere una loro evoluzione collegata all'età, in particolare, sembra che le risposte tendano a decrescere proporzionalmente all'aumentare della variabile *eta* e ciò è coerente con la natura dei dati. Infatti, è risaputo e comprovato da studi medici, che col passare degli anni la statura tenda a ridursi, e con essa la lunghezza delle ossa e di tutte le parti del corpo umano.

Inoltre, è ragionevole pensare che ci siano differenze tra i due sessi nelle grandezze osservate. Infine, nel caso di particolari patologie, si possono riscontrare variazioni significative delle dimensioni di grandezze corporee.

Si può, per altro, ritenere sensato che ci sia un effetto anche delle interazioni: in genere, alcune patologie sono più comuni in età avanzata. L'effetto congiunto di patologia ed età risulterà quindi rilevante; inoltre, se ci sono differenze tra i due sessi nelle misurazioni antropometriche, è logico anche pensare che tali discrepanze possano aumentare proporzionalmente all'età.

In conclusione, l'applicazione del modello di regressione multivariata ai dati porta a concludere che esiste una relazione significativa tra la misurazione delle grandezze antropometriche più rilevanti e l'età, il sesso e la presenza o meno di determinate patologie. È quindi possibile fornire stime delle grandezze che rappresentano le variabili risposta del caso-studio, applicando il modello ritenuto ottimale, a partire dai valori osservati delle esplicative. Inoltre, poiché l'interesse degli antropologi era rivolto soprattutto ai soggetti in età avanzata, è anche possibile, a partire dai dati e dal modello stimato, costruire curve di crescita, che mettano in evidenza l'evoluzione delle variabili risposta al variare dell'età.

Appendice

Distribuzioni multivariate

Distribuzione di Wishart

Sia \mathbf{M} una matrice ($p \times p$) che può essere scritta come $\mathbf{X}^T \mathbf{X}$, dove \mathbf{X} è una matrice ($m \times p$) che segue una distribuzione normale p -variata con media nulla e matrice di varianza e covarianza Σ . Allora \mathbf{M} ha distribuzione di Wishart con matrice di scala Σ e m gradi di libertà. Se $\Sigma = \mathbf{I}_p$, la distribuzione si dice standard.

Nel caso in cui p sia uguale ad 1, la distribuzione $W_1(\sigma^2, m)$ è data da $\mathbf{x}^T \mathbf{x}$, con \mathbf{x} vettore ($m \times 1$) distribuito come una da $N(0, \sigma^2)$. In questo caso $W_1(\sigma^2, m)$ è la distribuzione χ_m^2 moltiplicato per σ^2 .

La matrice di scala Σ gioca, per la distribuzione di Wishart, lo stesso ruolo che ha σ^2 nel caso particolare in cui $p = 1$.

Se $\mathbf{M} \sim W_p(\Sigma, m)$ e \mathbf{a} è un vettore ($p \times 1$) tale per cui $\mathbf{a}^T \Sigma \mathbf{a} \neq 0$, allora

$$\mathbf{a}^T \mathbf{M} \mathbf{a} / \mathbf{a}^T \Sigma \mathbf{a} \sim \chi_m^2.$$

Distribuzione A

Sia $\mathbf{A} \sim W_p(\mathbf{I}, m)$ e sia $\mathbf{B} \sim W_p(\mathbf{I}, n)$, indipendenti. Se $m \geq p$ allora

$$A = |\mathbf{A}| / |\mathbf{A} + \mathbf{B}|$$

ha distribuzione $A(p, m, n)$, dove m rappresenta il numero di gradi di libertà al numeratore ed n il numero di gradi di libertà al denominatore. La quantità $m + n$ rappresenta il numero di gradi di libertà totali. La distribuzione è invariante a cambiamenti di scala.

Un'ulteriore distribuzione legata alla A è la *Beta*. Infatti, se $\mathbf{u}_i \sim B\left(\frac{1}{2}(m + i - p), \frac{1}{2}p\right)$,

con $i = 1, \dots, n$, e se $\mathbf{u}_1, \dots, \mathbf{u}_n$ sono indipendenti, allora

$$\prod_{i=1}^n u_i \sim \Lambda(p, m, n).$$

La distribuzione *Beta* si ottiene dal rapporto di una variabile u e la somma tra u e v , con $u \sim \chi_p^2$ e $v \sim \chi_q^2$. In questo caso la *Beta* ha parametri $\frac{1}{2}p$ e $\frac{1}{2}q$.

Infine, la distribuzione Λ è anche legata alla F . Infatti

$$\frac{1 - \Lambda(p, m, 1)}{\Lambda(p, m, 1)} \sim \frac{p}{m - p + 1} F_{p, m-p+1}.$$

Distribuzione θ

Se $\mathbf{A} \sim W_p(\mathbf{I}, m)$ e $\mathbf{B} \sim W_p(\mathbf{I}, n)$, con \mathbf{A} e \mathbf{B} indipendenti ed $m \geq p$, allora il più grande degli autovalori di $(\mathbf{A} + \mathbf{B})^{-1}\mathbf{B}$ ha distribuzione $\theta(p, m, n)$, con m ed n che rappresentano, rispettivamente, il numero di gradi di libertà del numeratore e al denominatore.

Tale distribuzione ha alcune proprietà:

- $\theta(p, m, n)$ e $\theta(p, m + n - p, p)$ hanno la stessa distribuzione;
- $\frac{\theta(p, m, 1)}{1 - \theta(p, m, 1)} \sim \frac{1 - \Lambda(p, m, 1)}{\Lambda(p, m, 1)}$;
- $\frac{\theta(p, m, 1)}{1 - \theta(p, m, 1)} \sim \frac{p}{m - p + 1} F_{p, m-p+1}$.

Distribuzione di T^2 Hotelling

Se una statistica α può essere scritta come $m\mathbf{d}^T\mathbf{M}^{-1}\mathbf{d}$, dove \mathbf{d} e \mathbf{M} sono indipendenti e distribuite rispettivamente come una $N_p(0, \mathbf{I})$ e una $W_p(\mathbf{I}, m)$, allora α ha distribuzione T^2 di Hotelling di parametri p e m . Si può quindi scrivere $\alpha \sim T^2(p, m)$. Più in generale, se \mathbf{d} è una normale p -variata di media $\boldsymbol{\mu}$ e varianza $\boldsymbol{\Sigma}$ e \mathbf{M} proviene da una distribuzione Wishart con matrice di scala $\boldsymbol{\Sigma}$, allora si può scrivere

$$m(\mathbf{d} - \boldsymbol{\mu})^T\mathbf{M}^{-1}(\mathbf{d} - \boldsymbol{\mu}) \sim T^2(p, m)$$

Inoltre, per quanto riguarda le relazioni esistenti con altre distribuzioni notevoli, con $\alpha \sim T^2(p, m)$, si ha

$$\alpha = m \chi_p^2 / \chi_{m-p+1}^2 = \{mp/(m - p + 1)\} F_{p, m-p+1}.$$

Distribuzione t di Student multivariata

Sia \mathbf{X} distribuita come una normale p -variata con media $\boldsymbol{\mu}$ e matrice di varianza e covarianza $\boldsymbol{\Sigma}$ e indipendente da \mathbf{y} , distribuita come un χ_v^2 . Posto $\mathbf{u}_i = \mathbf{x}_i/(\mathbf{y}/v)^{1/2}$, per $i = 1, \dots, p$, allora il vettore \mathbf{u} segue una distribuzione t di Student multivariata con vettore delle medie $\boldsymbol{\mu}$, matrice di varianza e covarianza $\boldsymbol{\Sigma}$ e v gradi di libertà ossia

$$\mathbf{u} \sim t(v, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Per maggiori chiarimenti sulle distribuzioni che sono state qui brevemente introdotte, e su altre ad esse correlate, si veda Mardia *et al.*(1979).

Bibliografia

- [1] Anderson, T. W. (1958) *Introduction to multivariate statistical analysis*;
- [2] Box, G. E. P. (1949) *A general distribution theory for a class of likelihood criteria*;
- [3] Di Fonzo T. , Lisi F. (2005) *Serie Storiche Economiche*;
- [4] Levene, H. (1960). In *Contributions to Probability and Statistics: Essays in Honor of Harold Hotelling*, I. Olkin et al. eds., Stanford University Press;
- [5] Miles, D. B. (2003) *Multivariate Analysis of Variance*;
- [6] Rowe , D. B. *Multivariate Regression Generalized Likelihood Ratio Test for FMRI*;
- [7] Schatz, P.(2006) *Manova*;
- [8] Seber, G. A. F. (1984) *Multivariate observations*;
- [9] Vittadini, G. (1999). *On the use of Multivariate Regression models in the context of Multilevel Analysis*.