

Problema 1.

- a) Il modello ANOVA che andiamo a ipotizzare è un modello ANOVA con due fattori. Il primo è legato alla città e il secondo al tipo, il modello che creiamo è di conseguenza:
$$X_{ijk} = \mu + \tau_i + \beta_j + \gamma_{ij} + \epsilon_{ijk}$$

Dove μ indica la media totale τ l'effetto di essere in città diverse β l'effetto dovuto ai tipi diversi e infine γ rappresenta l'interazione tra i due fattori esposti sopra. Per assunzione ϵ devono essere normali indipendenti e identicamente distribuiti con media zero e varianza comune σ^2 .

Le assunzioni sono verificate attraverso uno Shapiro test nei singoli gruppi per la normalità e un Bartlett test per verificare l'uguaglianza delle varianze

- b) Come test per valutare la possibilità di rimuovere alcune covariate utilizziamo prima il test sulle interazioni quindi abbiamo come ipotesi nulla che $\gamma_{ij} = 0$ e come $H_1: \gamma_{ij} \neq 0$. E costruiamo il test basandoci sulla statistica test $SS_{\text{interaction}} / ((g-1) * (b-1)) / (SS_{\text{res}} / (gb(n-1)))$ se questa è maggiore di $F_{1-\alpha}((g-1) * (b-1), gb(n-1))$ allora l'ipotesi nulla è rifiutata. Nel nostro caso sia b che g sono pari a 2. Dal summary dell'ANOVA evinciamo che le interazioni possono essere trascurate e procediamo con un modello additivo a questo punto ci basiamo su una statistica test $SS_{\text{treatment}} / (g-1) / (SS_{\text{res}} + SS_{\text{interaction}} / (gb(n-b-g+1)))$ allo stesso modo rifiutiamo se questa è maggiore di $F_{1-\alpha}((g-1), gb(n-b-g+1))$. Dal summary dell'ANOVA possiamo considerare anche l'effetto del trattamento relativo alla città nullo quindi il nostro modello ridotto sarà del tipo
$$X_{ijk} = \mu + \beta_j + \epsilon_{ijk}$$
 dove le ϵ devono essere normali indipendenti e identicamente distribuiti con media zero e varianza comune σ^2 .
- c) Dobbiamo ora valutare gli intervalli di Bonferroni per la differenza delle medie in base al tipo di abito. In particolare dobbiamo fornire un solo intervallo di confidenza IC 0.95 $(\mu_1 - \mu_2) = [\text{media campionaria 1} - \text{media campionaria 2} \pm (t_{0.95}(n-2)) * \sqrt{(SS_{\text{res}}^2 / (n-2)) * (1/n_1 + 1/n_2)}]$ = [14.66669 15.15225]

Problema 2.

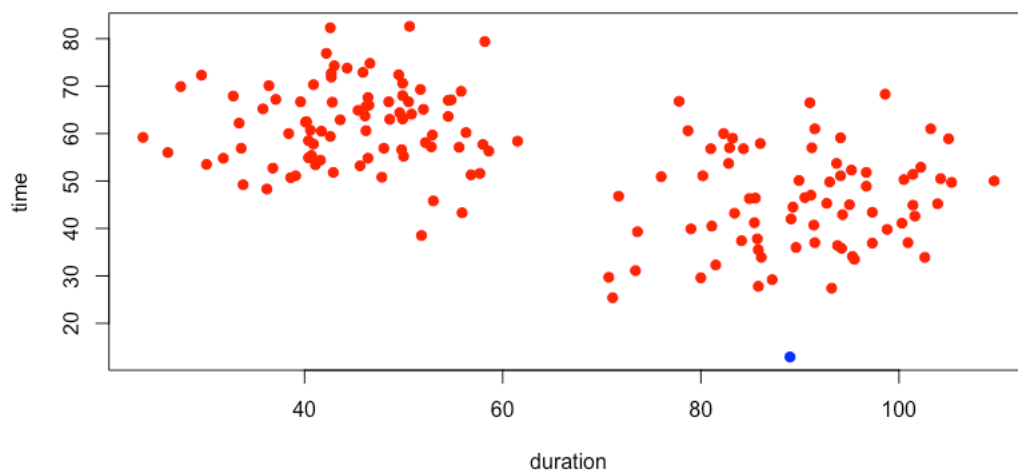
- a) Risolviamo il problema nel contesto dei Paired Data. In particolare abbiamo che le due misurazioni sono fatte sulle stessa unità statistica e quindi dobbiamo andare a valutare le differenze nelle due misurazioni. Definiamo $D_i = X_{i1} - X_{i2}$ ipotizziamo che le differenze siano normali e lo verifichiamo con un Shapiro test. A questo punto possiamo valutare la differenza tra le due medie nelle diverse variabili basandoci su un test per la media di una normale multivariata
 $H_0: \mu = 0$ vs $H_1: \mu \neq 0$
Lo basiamo sulla statistica test $m * (D_{\text{sample}})^2 / (S^2 * (D_{\text{sample}}))$ e valutiamo il p-value del nostro test. Il p-value dalla computazione è nullo di conseguenza rifiutiamo l'ipotesi nulla.
- b) Per avere i quattro intervalli di confidenza ci rifacciamo alla stessa quantità pivotale di prima in particolare avremo che
$$\text{SimCI}_{0.95}(\mu_i) = [(D_{\text{sample}})_i \pm \sqrt{((n-1)p / (n-p)) * F_{0.95}(p, n-p) * (S^2)_{ii} / n}]$$

	inf	center	sup
rice	-37.34390	-2.145937	33.05202
sashimi	107.51324	123.577812	139.64238
vegetables	-12.40493	11.463438	35.33181
okashi	111.98999	118.150000	124.31001

Possiamo vedere che sia l'intervallo di confidenza del riso che delle verdure contiene lo zero di conseguenza non abbiamo evidenza statistica per dire che il loro consumo sia diverso in diversi momenti. Mentre il consumo di sashimi e di okashi aumenta di molto nel periodo di festa.

Problema 3

- a) Vogliamo creare un clustering riguardante il successo o meno di un tour in base alle covariate e caratteristiche del tour che conosciamo. In particolare selezioniamo come distanza tra due unità statistiche la loro distanza euclidea e come distanza tra i cluster sfruttiamo il metodo single linkage una volta definite queste due quantità possiamo direttamente applicare un algoritmo di clustering gerarchico agglomerativo.



Where the centers are:

- For the red cluster: (66.07925, 54.33019)
- For the blue cluster: (89, 12.9)

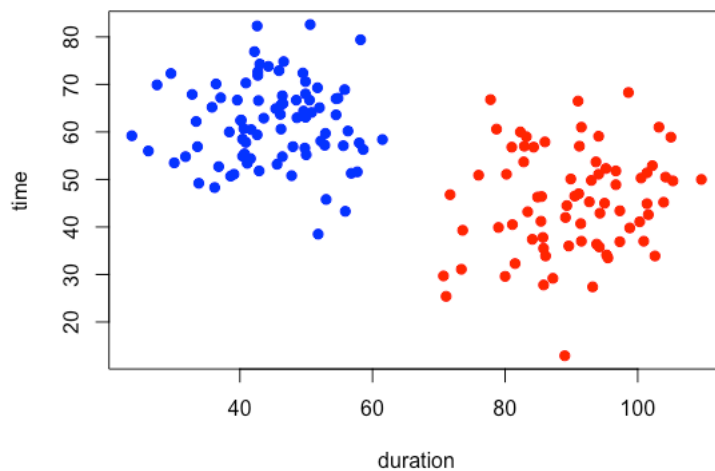
The numerosity of the clusters are:

- For the red cluster: 159
- For the blue cluster: 1

The cophenetic coefficient is:

0.8781562

- b) The clustering produced at point a is unsatisfactory since the single linkage is not good when the two clusters are not well separated we can help and have a better cluster with average or complete linkage. In this case :



Where the centers are:

- For the red cluster: (90.23867, 45.20000)
- For the blue cluster (45.03176, 61.89882)

The numerosity of the clusters are:

- For the red cluster: 75
- For the blue cluster: 85

The cophenetic coefficient is:

0.8807631

- c) Una volta divisi i dati nei due gruppi possiamo andare a controllare se sono normali e se hanno la stessa covarianza. In particolare grazie a un mc.shapiro test possiamo vedere che entrambi i gruppi sono normali ma non hanno la stessa covarianza quindi supponiamo che n sia abbastanza grande per valutare questi intervalli in modo asintotico

SimCI 0.95($\mu_1 - \mu_2$)i = [sample mean 1- sample mean 2 +- z 1-0.05/4 * sqrt(1/n1 * sigma 1 + 1/n2 * sigma 2)]ii]

BonfCI 0.95 (μ_1) = [sample mean 1+- t 1-0.05/4 (n1-1) * sqrt(1/n1 * sigma 1)]ii]

Diff.mean.duration 42.12192 48.29188

Diff.mean.starting time -20.21498 -13.18266

Mean.duration 87.83011 92.64723

Mean.Starting time 42.31408 48.08592

- d) Quello che possiamo notare dagli intervalli di confidenza è che tendenzialmente tour che durano di più e iniziano prima hanno maggior successo. Di conseguenza una strategia di successo potrebbe essere quella di valutare tour che partano alle 16:45 e durino 90 minuti.

Problema 4

- a) We based on the model

$$E = \beta_0 + \beta_1 \cdot x_1 + \beta_2 \cdot x_2 + \beta_3 \cdot x_3 + \beta_4 \cdot x_4 + \varepsilon,$$

with E the extension of the garden, x_1, x_2, x_3, x_4 the number of carps, maple trees, cherry trees and stones respectively, and $\varepsilon \sim N(0, \sigma^2)$.

In particular our design matrix is formed by $Z = [1 \ x_1 \ x_2 \ x_3 \ x_4]$ the estimation of beta is performed by OLS so that

$$\text{Beta} = (Z'Z)^{-1}Z'Y$$

Where Y is the vector of the observed response so E.

(Intercept)	carps	maple	cherry	stones
1442.86188	16.34107	28.69118	13.12381	14.02307

The residuals are normal and has a good cloud shape

- b) We need to perform a test of significance on linear combination of the beta and in particular two of them at the time so we perform to test like:

$H_0: C\beta = 0$ vs $H_1: C\beta \neq 0$

We use as test statistic $(C\beta)' (C(Z'Z)^{-1}C')^{-1} C\beta / p \cdot S^2$ and we compute the pvalue of the test.

If $C = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$

Then the p-value is really small so that we reject H_0 . The same hold for

$C = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}$

- c) The weakness of this model is collinearity indeed if we compute the vif of all the regressors we get
- | carps | maple | cherry | stones |
|----------|----------|----------|----------|
| 19.92811 | 18.87960 | 15.35421 | 17.45596 |

They are all greater than 10 so we have a lot of collinearity in the data. One approach to avoid this problematic is to use PCA regression

So we first find the PC that are orthogonal and then we base our analysis on the scores of the data.

Computing the PCA we can reduce the dimensionality to the first two PC. The first refer to a general weighed sum while the second to contrast between lake element and ground element.

Given b_0, b_1, b_2 form the PCA regression we need to recompute the value of the beta for the initial model

$$\begin{aligned} y &= b_0 + b_1 \cdot PC1 + b_2 \cdot PC2 + \epsilon = b_0 + b_1(e_{11}(X_1 - m_1) + e_{21}(X_2 - m_2) + e_{31}(X_3 - m_3) + e_{41}(X_4 - m_4)) + b_2(e_{12}(X_1 - m_1) + e_{22}(X_2 - m_2) + e_{32}(X_3 - m_3) + e_{42}(X_4 - m_4)) + \epsilon \\ &= b_0 - b_1 \cdot e_{11} \cdot m_1 - b_2 \cdot e_{12} \cdot m_1 - b_1 \cdot e_{21} \cdot m_2 - b_2 \cdot e_{22} \cdot m_2 - b_1 \cdot e_{31} \cdot m_3 - b_2 \cdot e_{32} \cdot m_3 - b_1 \cdot e_{41} \cdot m_4 - b_2 \cdot e_{42} \cdot m_4 + (b_1 \cdot e_{11} + b_2 \cdot e_{12}) \cdot X_1 + (b_1 \cdot e_{21} + b_2 \cdot e_{22}) \cdot X_2 + (b_1 \cdot e_{31} + b_2 \cdot e_{32}) \cdot X_3 + (b_1 \cdot e_{41} + b_2 \cdot e_{42}) \cdot X_4 + \epsilon \end{aligned}$$

(Intercept)

1484.613

Carps

15.64089

Maple

19.25401

Cherry

21.3864

Stones

15.08058