

## Statistica Computazionale Progredito esame del 3 febbraio 2020

**Istruzioni:** lo studente deve produrre un file in formato pdf in cui riporta il nome, cognome, numero di matricola e, per ogni esercizio a cui risponde, il codice R utilizzato per produrre il risultato, il risultato (valori, grafici,...), i commenti e quant'altro richiesto dall'esercizio.

1. I dati nell'oggetto `Clotting` si riferiscono ad un esperimento sul tempo di coagulazione del sangue. In particolare, è stato misurato il tempo di coagulazione del plasma sanguigno (in secondi) (`tempo`) in 18 campioni di plasma normale diluito con plasma privo di protrombina in modo da ottenere 9 diverse concentrazioni percentuali (`u`). La coagulazione è stata indotta con due diversi lotti di tromboplastina (`lotto`).
  - (a) Stimare un modello di regressione Gamma con link canonico (`inverse`) in cui `tempo` è la variabile risposta e `u` e `lotto` sono le variabili esplicative. Commentare il modello stimato e valutare l'eventuale inserimento di un termine di interazione tra le variabili esplicative.
  - (b) Effettuare un bootstrap delle unità statistiche e commentare i risultati ottenuti, in particolare con riferimento alla variabile `lotto` di cui si vuole calcolare anche un intervallo di confidenza di livello 0.95. Si commenti anche la scelta del tipo di bootstrap rispetto alla natura dei dati.
  - (c) Il modello stimato nei punti precedenti assume una funzione legame inversa tra la media del tempo di coagulazione e il predittore lineare. Si vuole valutare una funzione di legame alternativa, in particolare si assume una funzione legame logaritmica (`log`). Per confrontare i due modelli si può utilizzare come test la differenza delle log verosimiglianze stimate. Dare il valore del test osservato nel campione. (Suggerimento: la funzione `logLik` potrebbe essere utile...).
  - (d) Trovare la distribuzione nulla stimata simulata del test al punto precedente attraverso un bootstrap parametrico, quando si vuole verificare l'ipotesi nulla che il modello abbia legame inverso contro l'alternativa che il legame sia logartmico. Dare una stima del relativo  $p$ -value. (Suggerimento: se `modello` è il modello stimato, `modello$fitted.values` restituisce le medie stimate della variabile risposta,  $\hat{\mu}_i$ , e `summary(modello)$dispersion` è la stima del parametro di dispersione,  $\tilde{\phi}$ ; inoltre  $E(Y_i) = \mu_i$  e  $\text{Var}(Y_i) = \phi\mu_i^2$ )
  - (e) Ripetere il punto precedente invertendo il ruolo delle due ipotesi.
  - (f) Alla luce dei  $p$ -value calcolati ai due punti precedenti, quale funzione legame sembra preferibile?

2. L'oggetto **Venicemax** contiene dati relativi ai livelli massimi annui della marea registrati a Venezia tra il 1875 e il 2019. In particolare, la variabile **valore** indica il valore massimo della marea (in cm) e la variabile **anno** il corrispondente anno. Si assume che i massimi annuali siano un campione casuale semplice da una distribuzione Gumbel con funzione di ripartizione

$$F_Y(y; \mu, \sigma) = \exp[-\exp\{-(y - \mu)/\sigma\}],$$

con  $\mu \in \mathbb{R}$ ,  $\sigma > 0$  e  $y \in \mathbb{R}$ .

- (a) Mostrare che la log verosimiglianza per  $\theta = (\mu, \sigma)$  è pari a

$$\ell(\theta) = \begin{cases} -n \log \sigma - \sum_{i=1}^n \frac{y_i - \mu}{\sigma} - \sum_{i=1}^n \exp\left\{-\frac{y_i - \mu}{\sigma}\right\}, & \text{se } \sigma > 0, \\ -\infty, & \text{altrimenti.} \end{cases}$$

Scrivere una funzione in R che calcola  $\ell(\theta)$ .

- (b) Trovare numericamente la stima di massima verosimiglianza di  $\theta$  e dare una valutazione numerica dello standard error di ciascuna componente della stima.  
(c) Fare un grafico della log verosimiglianza in un'opportuna regione dello spazio parametrico.  
(d) Indicato con  $y_p$  il quantile di  $Y$  che lascia probabilità  $p$  a destra, dimostrare che è pari a

$$y_p = \mu - \sigma \log\{-\log(1 - p)\}.$$

Trovare la stima di massima verosimiglianza di  $y_p$ , quando  $p = 0.001$ .

- (e) Scrivere una funzione che calcola la log verosimiglianza profilo per  $y_{0.001}$  e farne un grafico in un intervallo opportuno. Trovare anche un intervallo di confidenza di livello 0.95 per  $y_{0.001}$  basato sul log rapporto di verosimiglianza.  
(f) Si assumono per  $\mu$  e  $\sigma$  distribuzioni a priori indipendenti con

$$\mu \sim U(a_1, b_1), \quad \sigma \sim \text{Gamma}(a_2, b_2),$$

con  $a_1 = 80$ ,  $b_1 = 250$ ,  $a_2 = b_2 = 0.01$ . Dare un'interpretazione alla scelta della distribuzione a priori. Sembra una distribuzione particolarmente non informativa?

- (g) Scrivere una funzione che calcola la distribuzione a posteriori per  $\theta$  ed implementare un algoritmo Metropolis-Hastings per simulare da tale distribuzione. Scegliere i parametri dell'algoritmo in modo da ottenere una convergenza soddisfacente. Controllare quest'ultima sia con indicatori numerici che grafici, possibilmente usando più catene e commentando.  
(h) Utilizzare i valori simulati al punto precedente per dare una
- i. rappresentazione grafica della distribuzione a posteriori di  $\theta$ ,
  - ii. stima delle medie e degli standard error a posteriori di  $\mu$  e  $\sigma$ ,
  - iii. stima di intervalli di credibilità di probabilità 0.95 per  $\mu$  e  $\sigma$ ,
  - iv. stima della media, di un intervallo di credibilità e una rappresentazione grafica della distribuzione a posteriori di  $y_{0.001}$ . Fare un confronto con i risultati frequentisti ottenuti in precedenza e commentare.
  - v. Dare un intervallo di previsione *equi-tailed* di probabilità 0.95 per il valore massimo della marea nel 2020. (Suggerimento: è facile simulare da  $Y$  tramite inversione),