

Statystyka II Mat. Egzamin komputerowy 13.06.08

1. Wczytaj z archiwum *MacierzePot4.zip* 12 potencjałów kontaktowych (PK). PK są to symetryczne macierze 20×20 opisujące oddziaływania 210 par aminokwasów w białkach. Napisz ogólną funkcję `lowerVec`, która przekształca macierz na wektor złożony z elementów należących do jej dolnej części bez przekątnej. Ze zbioru wektorów $\mathbf{x}_1, \dots, \mathbf{x}_{12}$ otrzymanych z 12 PK, zbuduj macierz danych $X = [\mathbf{x}_1, \dots, \mathbf{x}_{12}]$ (190 obserwacji 12-wymiarowych). Za pomocą funkcji `boxplot.stat` sprawdź, czy w zbiorze X są obserwacje odstające?
2. Wykonaj klasteryzację hierarchiczną `hclust` metodą "centroid" danych z poprzedniego zadania. Dendrogram zapisz do pliku `.eps`. Na podstawie wykresu separowalności i sylwetki oszacuj liczbę klastów. Napisz ogólną funkcję `reduceDim(X, proc)`, która przekształca macierz danych X na macierz złożoną z jej k -pierwszych składowych głównych, gdzie k jest minimalną liczbą składowych, których łączna wariancja stanowi przynajmniej `proc` procent całkowitej wariancji X . Jakie jest k dla `proc=80` ?
3. Porównaj metody klasyfikacji `lda`, `qda`, `glm` (regresja logistyczna) w eksperymencie krosvalidacji pięciokrotnej na danych `MathAchSchool[, -c(1, 6)]` z biblioteki `nlme`. Cecha przewidywana y przyjmuje dwie wartości: „Catholic” i „Public”. Jako miarę efektywności predykcji przyjmij współczynnik Goodmana-Kruskala między y oraz y_{pred} .
4. Wykorzystując rozkład $X=QR$, bez użycia pętli, bez mnożenia i odwracania macierzy, napisz własną wersję funkcji `hatvalues`. Argumentem funkcji ma być tylko macierz cech obserwowalnych (predyktorów) X .
5. Zbuduj model regresji liniowej wielu zmiennych dla przefiltrowanych danych `Cars93`:
`library(MASS) ; X=na.omit(Cars93[,c(2,7,12:15,17,19,21,22,24,25)]);`
`row.names(X)=X[,1]; X=X[,-1]`. Cechą przewidywaną jest `MPG.city`.