

8.1 Dane w pliku *uscrime.txt* zawierają informacje dotyczące 47 stanów w USA.

- (a) Dopasuj model L0 opisujący zależność liniową pomiędzy zmienną R (wskaźnik przestępczości) a pozostałymi zmiennymi. Pomiń zmienną Ex0 (jest ona silnie skorelowana ze zmienną Ex1).
- (b) Dopasuj model L1 w którym są zmienne: Age, Ed, Ex1, U2, W, X.
- (c) Dopasuj model L2 w którym zostaną tylko dwie zmienne objaśniające Ed i Ex1.
- (d) Oszacuj błąd predykcji metodą krosvalidacji typu leave-one-out dla modelu L0, L1, L2. W którym modelu mamy najmniejszy a w którym największy błąd predykcji?

8.2 Wczytaj dane *longley*. Zbiór zawiera informacje o liczbie osób zatrudnionych w USA w latach 1947-1962.

- (a) Dopasuj model liniowy opisujący zależność liniową między zmienną Employed a pozostałymi zmiennymi.
- (b) Oblicz estymatory parametrów używając metody regresji grzbietowej korzystając z definicji oraz rozkładu SVD.
- (c) Dopasuj model regresji grzbietowej korzystając z pakietu *glmnet*. Wybierz optymalną wartość parametru λ stosując metodę krosvalidacji. Zbadaj stabilność estymacji dla różnych wartości parametru λ . Oblicz estymatory współczynników dla $\lambda = 0.03$.
- (d) Porównaj współczynnik przy zmiennej GNP (Produkt narodowy brutto) dla modelu z punktu (a) oraz dla modelu dopasowanego przy pomocy regresji grzbietowej z parametrem $\lambda = 0.03$.

8.3 Wczytaj dane z pliku *prostate.txt*. Dane zawierają informacje dotyczące raka prostaty u 97 mężczyzn.

- (a) Dopasuj model liniowy opisujący zależność zmiennej lpsa (logarithm of prostate specific antigen) od pozostałych zmiennych (z wyjątkiem zmiennej train). Dokonaj selekcji zmiennych metodą eliminacji "wstecz", z kryterium AIC.
- (b) Dopasuj model używając metody lasso. Przeanalizuj zachowanie estymatorów w zależności od parametru λ .
- (c) Na podstawie metody lasso i krosvalidacji, dokonaj selekcji zmiennych. Wyświetl które zmienne są po kolei dołączane do modelu gdy rośnie wartość parametru λ .
- (d) Oblicz wartości estymatorów w wybranym modelu.

8.4 Udowodnij że:

- (a) Estymator $\hat{\beta}^{RIDGE}$ jest obciążony.
- (b) Estymator RIDGE ma mniejszą wariancję niż estymator MNK, tzn.:

$$var(\hat{\beta}^{RIDGE}) \leq var(\hat{\beta}^{MNK}).$$

8.5 Zaproponuj eksperyment symulacyjny, który ilustruje powyższy wynik teoretyczny.