

Projekt Zaliczeniowy 1

Ania Macioszek, Dorota Celińska-Kopczyńska

Dane do projektu są w pliku `people.tab`. Są to dane symulowane; opisują wiek, wagę, wzrost, płeć, stan cywilny, liczbę dzieci, posiadane zwierzę domowe oraz miesięczne wydatki pewnych osób.

1. Wczytaj dane, obejrzyj je i podsumuj w dwóch-trzech zdaniach. Pytania pomocnicze: ile jest obserwacji, ile zmiennych ilościowych, a ile jakościowych? Czy są zależności w zmiennych objaśniających (policz i podaj VIF)? Czy występują jakieś braki danych? **(1 pkt)**

2. Podsumuj dane przynajmniej trzema różnymi wykresami. Należy przygotować:

- wykres typu scatterplot dla wszystkich zmiennych objaśniających ilościowych i zmiennej objaśnianej
- Wykresy typu pudełkowy (boxplot) dla wybranej zmiennej ilościowej
- Wykres typu słupkowy (barplot) dla wybranej zmiennej jakościowej

Mile widziane dodatkowe wykresy wg własnej inwencji (np histogram, punktowy, liniowy, mapa ciepła...). **(3 pkt)**

3. Podaj przedziały ufności dla wartości średniej i wariancji dla zmiennych wiek i wzrost. Jeżeli w celu wyliczenia przedziału ufności musisz poczynić jakieś założenia (np. założyć że zmienna pochodzi z rozkładu normalnego), zaznacz to i skomentuj czy wydaje Ci się to w danym przypadku uprawnione. Opisz wszelkie dodatkowe operacje, jakie zostały wykonane przed testem (takie jak usunięcie obserwacji odstających). Przedyskutuj, dla której ze zmiennych oczekujesz prawidłowych wyników. **(1 pkt)**

4. Sformułuj i zweryfikuj cztery hipotezy:

- dotyczącą różnicy między średnią wartością wybranej zmiennej dla kobiet i dla mężczyzn
- dot. niezależności między dwoma zmiennymi ilościowymi
- jedną dot. niezależności między dwoma zmiennymi jakościowymi
- jedną dot. rozkładu zmiennej (np. "zmienna A ma rozkład wykładniczy z parametrem 10")

Każda hipoteza po **2 punkty** (w sumie **8 pkt**). Punktowane jest sformułowanie hipotezy zerowej i alternatywnej, wybranie właściwego testu, przeprowadzenie testu i podjęcie decyzji czy odrzucamy hipotezę zerową.

4. Oszacuj model regresji liniowej, przyjmując za zmienną zależną (y) wydatki domowe (expenses) a zmienne niezależne (x) wybierając spośród pozostałych zmiennych. Rozważ, czy

konieczne są transformacje zmiennych lub zmiennej objaśnianej. Podaj RSS , R^2 , p-wartości i oszacowania współczynników i wybierz właściwe zmienne objaśniające, które najlepiej tłumaczą Household_expenses. Sprawdź czy w wybranym przez Ciebie modelu spełnione są założenia modelu liniowego i przedstaw na wykresach diagnostycznych: wykresie zależności reszt od zmiennej objaśnianej, na wykresie reszt studentyzowanych i na wykresie dźwigni i przedyskutuj, czy są spełnione. **(2 pkt)**.

Wynikiem ma być raport w formacie .Rmd oraz skompilowany do html.