

# Dokumentacja projektu - PADR

Joanna Rancew

## Heart Disease Dataset

Dane analizowane w projekcie pochodzą z ogólnodostępnej bazy danych i dotyczą aspektów związanych z chorobą serca.

## Cele projektu

Celem projektu była analiza oraz prezentacja danych przedstawionych w zbiorze “Heart Disease Dataset”. Podczas pracy nad programem wielokrotnie miałam okazję nabyć nowe informacje i wykorzystać zdobytą wiedzę w praktyce.

**Z założenia program miał spełniać poniższe funkcje:**

- poprawne wczytywanie i obróbka danych
- połączenie zbiorów danych
- wyświetlenie danych w postaci histogramu w zależności od wybranych parametrów
- przewidywanie diagnozy z użyciem drzew decyzyjnych
- prezentacja danych w postaci boxplotu
- prezentacja tekstowa przewidywanej diagnozy po wypełnieniu testu
- korelogram - prezentacja korelacji wyników badań i diagnozy
- zastosowanie możliwości xgBoost do analizy i klasyfikacji

## Źródła danych

- Hungarian Institute of Cardiology. Budapest: Andras Janosi, M.D.
- University Hospital, Zurich, Switzerland: William Steinbrunn, M.D.
- University Hospital, Basel, Switzerland: Matthias Pfisterer, M.D.

*Dane zostały odpowiednio dostosowane do użytku w programie*

## Informacje podstawowe

Baza danych zawiera 76 zmiennych. W oparciu o nią przeprowadzane były eksperymenty. Do eksperymentów brane były pod uwagę poniższe zmienne z bazy danych.

## Zmienne

1. age - wiek w liczbach
2. sex - plec, 1 - mezczyzna, 0 - kobieta
3. cp - typ bolu w klatce (1- typowa dusznosc, 2 - nietypowa, 3 - inny rodzaj bolu, 4 - bez objawow)
4. trestbps - cisnienie spoczynkowe krwi w mm Hg
5. chol - poziom cholesterolu w mg/dl
6. fbs - cukier we krwi (1- powyzej 120 mg/dl, 0 - ponizej 120 mg/dl)
7. restecg - 0 w normie, 1 - odchyłki od normalnosci na odcinku ST-T, 2 - przerost lewej komory wg kryteriow Estes
8. thalach - maksymalne tetno podczas wysilku
9. exang - czy wystapila dusznosc, dławica wysilkowa
10. diag - diagnoza (0- zdrowy, 1,2,3,4- rozny stopien zaawansowania choroby)
11. painloc - lokalizacja bolu (podsrodkowa - 1, inna - 0)
12. painexer - powod bolu (prowokowany wysilkiem - 1, inny - 0)
13. relrest - czy bol ustaje po odpoczynku
14. smoke - czy pacjent jest palacy
15. cigs - liczba papierosow dziennie
16. years - liczba lat palenia

Efekt pracy, opisywane w dokumentacji wykresy i program dostepne sa pod adresem: Moja aplikacja shiny

## Stosowane biblioteki

```
library(shiny)
library(dplyr)
library(ggplot2)
library(readxl)
library("stringr")
library(DT)
library(ggpubr)
library(tidyr)
library(corrplot)
library(rpart)
library(rpart.plot)
library(ggcorrplot)
library(party)
library(tree)
library(rattle)
```

## Działanie programu

### Wprowadzanie danych

Wprowadzanie danych realizowane jest w kilku etapach. Dane początkowo były przechowywane w nieuporządkowanym pliku `.data`, który wymaga odpowiedniego przygotowania. Na początku czytane są kolejne linie i segregowane po zmiennej `"name"`, która jest ostatnio zmienna każdego wiersza w oryginale.

```
data <- readLines("data/data_hungarian.data") #wczytuje linjki z pliku po kolei
#wylapuje indeksy kiedy pojawia sie "name"
ind <- grep("name",data)
data_proc <- vector("character", length(ind))
for(i in ind){

  data_proc[i/10] <- paste(data[(i-9)],data[i-8],data[i-7],data[i-6],data[i-5],data[i-4],
                           data[i-3],data[i-2],data[i-1],data[i])

  data_proc[i/10]<-gsub("name","NA",data_proc[i/10])
  data_proc[i/10]<-gsub("-9. ", "NA ",data_proc[i/10])
  data_proc[i/10]<-gsub("-9 ", "NA ",data_proc[i/10])
  data_proc[i/10]<-gsub(" ", ",",data_proc[i/10])
  #-9 oraz 9 oznaczaly nieprawidlowe lub brakujace wartosci
}
```

### Odpowiednie przekształcanie danych

Otrzymane dane wymagały przekształcenia, czyli:

- rozdzielenia elementów
- odpowiedniego zastąpienia wartości `NA`
- zmiany na typ `numeric`

```

df <- matrix(nrow=length(ind),ncol=76)
ind <- ind/10
for(i in ind){
  vector_test<-str_split(data_proc[i], ",")
  vector_test<- (unlist(vector_test))
  vector_test <- replace(vector_test, vector_test == "NA", NA)
  vector_test <- replace(vector_test, vector_test == "NA ", NA)
  vector_test<- as.numeric((vector_test))
  df[i,] <-vector_test
}

```

## Przekształcenie na ramke danych i nazwanie kolumn

```

df <- as.data.frame(df)
df <- cbind(df$V3,df$V4,df$V9,df$V10,df$V12,df$V16,
           df$V19,df$V32,df$V38,df$V40,df$V41,df$V51,
           df$V58, df$V5,df$V6,df$V7,df$V13,df$V14,
           df$V15,df$V17,df$V18)

Hungarian_data <- as.data.frame(df)
colnames(Hungarian_data) <- c("age","sex","cp","trestbps",
                              "chol","fbs","restecg","thalach","exang",
                              "oldpeak","slope","thal","diag","painloc",
                              "painexer","relrest","smoke","cigs","years",
                              "dm","famhist")

```

Dokładnie takie same kroki zostały wykonane dla plików *data\_va.data* oraz *data\_switzerland.data*. Ze względu na niedużą użyteczność dane z Cleveland zostały odrzucone na etapie wprowadzenia.

## Dostosowanie danych

Po wstępnym przetworzeniu danych, wymagały one dodatkowej obrobki. Na tym etapie dołączyłam dodatkową kolumnę *diag2*, która określała konkretnie, czy dany pacjent jest zdrowy chory (0-zdrowy, 1-chory). Zastąpiłam też niektóre dane brakujące średnią wartością pozostałych danych, co zaprezentowane zostało poniżej:

```
heart_data$diag2 <- heart_data$diag
heart_data$diag2[heart_data$diag2>0] = rep(1,length(heart_data$diag2[heart_data$diag2>0]))

heart_data$chol[heart_data$chol==0] <- mean(heart_data$chol, na.rm=TRUE)
heart_data$chol[is.na(heart_data$chol)] <- mean(heart_data$chol, na.rm=TRUE)

heart_data$trestbps[heart_data$trestbps==0] <- mean(heart_data$trestbps, na.rm=TRUE)

heart_data$fbs[is.na(heart_data$fbs)] <- rep(0,
      length(heart_data$fbs[is.na(heart_data$fbs)]))
heart_data$exang[is.na(heart_data$exang)] <- rep(0,
      length(heart_data$exang[is.na(heart_data$exang)]))

heart_data$relrest[is.na(heart_data$relrest)] <- rep(0,
      length(heart_data$relrest[is.na(heart_data$relrest)]))

heart_data = subset(heart_data, select = -c(slope,oldpeak,thal,dm,famhist))
```

## Sekcja serwera

Po wprowadzeniu danych, zajęłam się odpowiednią obsługą ich w aplikacji. Korzystałam z dwóch plików: *ui.R* oraz *server.R*, co pozwoliło na czytelne rozdzielenie interfejsu graficznego od działania serwera.

```
result <- heart_data
#modyfikacja danych obsługiwanych
dane <- reactive(dplyr::filter(result, sex == input$jaka_plec &
      age >= input$jaki_wiek[1] & age <= input$jaki_wiek[2]))
dane_palacz <- reactive(dplyr::filter(result, sex == 1 & smoke == input$czy_pali &
      age >= input$jaki_wiek[1] & age <= input$jaki_wiek[2]))
#gdy jest jako reactive - to traktowane jest jako funkcja
```

## Wykresy

W aplikacji, która odpowiada za obróbkę, analizę i prezentację danych, kluczowym elementem są wykresy. Zgodnie z założeniami projektu, prezentowałam dane o pacjentach w różny sposób. Miedzy innymi używając *boxplota*, *histogramu* i *wykresu gestosci wystepowania danej wartosci*.

**Prezentacja danych dotyczaca palacych lub niepalacych mezczyzn:** Dane zaprezentowały zostały przy użyciu funkcji *geom\_boxplot*, która pozwoliła dobrze zaprezentować różnice w spoczynkowym ciśnieniu krwi zdrowych i chorych, w zależności od tego, czy są palaczami czy nie.

```
ggplot(dane_palacz(),aes(y=trestbps, x=as.factor(diag2)))+ geom_boxplot(na.rm=TRUE)+  
  scale_x_discrete(labels=c("zdrowy", "chory"))+  
  labs(y="Spoczynkowe ciśnienie krwi",x="Występowanie choroby")+  
  theme(legend.position="top")+theme(legend.title = element_blank())
```

**Prezentacja danych dotyczących poziomu cholesterolu** Dane zaprezentowane zostały z użyciem funkcji *geom\_density*, która podobnie jak histogram, pozwala rzetelnie zaprezentować częstość występowania danych wartości.

```
ggplot(dane(),aes( x=chol, color=as.factor(diag2)))+ geom_density(na.rm = TRUE)+  
  scale_color_discrete(labels=c("zdrowy", "chory"))+  
  labs(x="poziom cholesterolu",y="gestosc wystepowania")+  
  theme(legend.position="top")+theme(legend.title = element_blank())
```

**Prezentacja danych dotyczących tetna - histogram** Dane dotyczące tetna zdrowych i chorych pacjentów, zostały zaprezentowane dzięki funkcji *geom\_histogram*, która pozwoliła zaprezentować zależność tego od stanu zdrowia pacjenta.

## Korelacja

Przed analizą i wykorzystaniem posiadanego zbioru danych do przewidywania diagnozy, kluczową kwestią jest ocena korelacji danych. W kolumnie *diag2* oraz w kolumnie *diag* umieszczone zostały diagnozy. Korelacja innych danych z nimi, pozwala ocenić, jak wyniki badań i styl życia wpływają na ryzyko choroby. Stosowane funkcje *cor* oraz *ggcorrplot* pozwoliły w przejrzysty sposób zaprezentować korelacje między danymi.

```
corr <- round(cor(result),1)  
ggcorrplot(corr)
```

## Sekcja interfejsu

Zaprezentowane powyżej fragmenty programu pozwalają wykonać działania serwera i obróbkę danych. Za prezentację danych w aplikacji shiny odpowiada kod programu *ui.R*. Składają się na niego poniższe elementy:

### Layout

Za odpowiednie ułożenie elementów w aplikacji odpowiada *fluidPage* oraz *sidebarLayout*.

### Wprowadzanie danych:

Dane stosowane do prezentacji i analizy wyników pacjenta wprowadzane są różnymi metodami. Poniżej przedstawie ich działanie:

```
sliderInput("jaki_wiek", "Wiek:",  
            min = 28, max = 77,  
            value = c(30,75)),
```

### SliderInput

```
selectInput("jaka_plec", "Plec", choices = c("Kobieta"=0, "Mezczyzna"=1), multiple=FALSE)
```

### selectInput

**Dane do testu** Dane do testu wprowadzane są przy użyciu *radioButtons* i *numericInput*.

```
radioButtons(  
  "test_sex",  
  "Plec",  
  inline = TRUE,  
  choiceNames = c("Mezczyzna", "Kobieta"),  
  choiceValues = c("1", "0")  
) ,  
numericInput(  
  "test_age",  
  "Wiek",  
  min=0,  
  value=40  
) ,
```

## Prezentacja danych

Prezentacja danych odbywa się w panelu głównym przy użyciu *tabsetPanel* oraz *tabPanel*, co porządkuje widoki.

```
tabPanel("Cholesterol a stan zdrowia", plotOutput("plot_density")),
  tabPanel("Tetno a stan zdrowia",
    plotOutput("plot_hist_chory"),
    plotOutput("plot_hist_zdrowy")
  ),

  tabPanel("Tetno - dla mezczyzn palacych i niepalacych", selectInput("czy_pali", "Wybierz c",
    plotOutput("distPlot"))
```

Diagnoza prognozowana po wypełnieniu testu wyświetlana jest jako tekst przy użyciu *textOutput*.

```
textOutput("text"),
```



## Skrypty analityczne - analiza i przewidywanie danych

Do analizy danych wykorzystywane były *drzewa decyzyjne*. Drzewo tworzone było za pomocą funkcji *rpart* na danych treningowych *training*. Następnie testowane na danych testowych *test*. Wyniki nie były satysfakcjonujące, lecz po przetestowaniu kilku metod, zdecydowałam się na pozostanie przy tej funkcji.

W trakcie tworzenia projektu badałam także rozwiązanie z wykorzystaniem funkcji *tree* oraz *ctree* po wyświetleniu schematu drzew i wyników, mogłam stwierdzić, że stosowanie *rpart* dla mojego zbioru danych jest wystarczające. Już sam korelogram pozwolił stwierdzić, że dane nie są od siebie silnie zależne i postawienie diagnozy może być utrudnione.

Podczas pracy nad dodatkowym zagadnieniem (drzewa decyzyjne, *xgBoost*) korzystałam z zasobów forum internetowych na temat programowania w R, prac na temat drzew decyzyjnych oraz dokumentacji. Zrezygnowałam z wykorzystania *xgBoost* w realizowanym projekcie ze względu na ograniczony czas i złożoność programu.

```
training <- heart_data[1:300,]
test <- heart_data[301:617,]
test$diag2 <- as.factor(test$diag2)
training$diag2 <- as.factor(training$diag2)
tree <- rpart(diag2~., training)
predicted_Classes <- predict(tree, test, type="class")
```

Dla danych pobranych w ramach wykonywanego testu (*our\_test*), przeprowadziliśmy weryfikację i prognozę diagnozy:

```
our_test$age=input$test_age
our_test$sex=as.integer(input$test_sex)
our_test$cp=as.integer(input$test_cp)
our_test$painloc=as.integer(input$test_painloc)
our_test$trestbps=input$test_trestbps
our_test$chol=input$test_chol
our_test$fbs=as.integer(input$test_fbs)
our_test$thalach=input$test_thalach
our_test$exang=as.integer(input$test_exang)
our_test$painexer=as.integer(input$test_painexer)
our_test$relrest=as.integer(input$test_relrest)

our_result <- predict(tree, our_test, type = "class")

if(our_result[1]==0)
  diagnosis = "Wynik: Masz mniejsze ryzyko zachorowania"
else if(our_result[1]==1)
  diagnosis="Wynik: masz duże ryzyko zachorowania"
else diagnosis="Zapraszam do wypełnienia testu"

#diagnosis - wyświetlane jako tekst w aplikacji
```

Diagnoza stawiana jest poprawnie, jednak samo wyświetlanie jej nie działa poprawnie.

## Podsumowanie

Praca nad projektem podzielona była na różne etapy. Najtrudniejszym z nich był, ku mojemu zaskoczeniu, etap wprowadzania danych, które były zapisane w niezbyt wygodny sposób w pliku. Po wprowadzeniu danych, przeanalizowałam je i odeszłam nieco od schematu proponowanego na stronie - sama oceniając, które dane są w moim programie potrzebne i użyteczne.

Następnym trudnym etapem był wybór danych do prezentacji i analizy. Sam korelogram pokazuje, że dane nie są od siebie mocno zależne, więc histogram czy inne wizualizacje nie do końca wskazywały na to, jakie czynniki wpływają na chorobę.

Ostatni etap - etap klasyfikacji i wykorzystanie drzew decyzyjnych, pozwolił mi poznać metody i rozszerzyć wiedzę, która jest bardzo dobrym wstępem do dalszego rozwoju w pisaniu skryptów analitycznych.