

Błażej Kucman
ind. 238228
gr.1

Inteligencja Obliczeniowa – Projekt 2

Zgłębianie danych

Tematem zadania było przeciwiczenie poznanych na laboratoriach technik zgłębiania danych na wybranym zbiorze danych (najlepiej z kolumną „class” z dwoma możliwymi wartościami).

1. Wstęp i omówienie bazy

Do zadania została wybrana baza z wyniki i danymi osób z problemami chorób serca. Na ich podstawie chcemy się dowiedzieć przy jakich wynikach oraz danych osoby jest możliwa choroba serca.

Baza składa się z kolumn:

Skrót ST – częstoskurcz zatokowy

age - wiek w latach,

sex – płeć(1=mężczyzna, 0= kobieta),

cp – rodzaj bólu w klatce piersiowej (wartości 1-4)

1. Typowa angina
2. Nietypowa angina
3. Nie anginowy ból
4. Bezobjawowy

trestbps - spoczynkowe ciśnienie krwi (w mm Hg przy przyjęciu do szpitala)

chol - poziom cholesterolu w surowicy w mg/dl

fbs – cukier we krwi(na czczo) > 120 mg/dl (1-tak, 0-nie)

restecg - spoczynkowe wyniki elektrokardiograficzne(wartości 0-3)

- Wartość 0: normalna
- Wartość 1: nieprawidłowość fali ST-T (odwrócenie załamka T i / lub uniesienie odcinka ST> 0,05 mV)
- Wartość 2: wykazanie prawdopodobnego lub określonego przerostu lewej komory według kryteriów Estes.

thalach – osiągnięte maksymalne tętno

exang – angina wywołana wysiłkiem (1-tak , 0-nie)

oldpeak – spadek ST wywołany pod czas ćwiczeń względem stanu spoczynku

slope – nachylenie spadku ST pod czas ćwiczeń(0-3)

- Value 1: upsloping
- Value 2: flat
- Value 3: downsloping

ca - liczba głównych naczyń (wartości 0-3) zabarwionych za pomocą fluoroskopii
thal – 3 – normalny; 6 – wyleczony , 7- wada odwracalna
target – 0-zdrowy, (1-4) – chory

Jako klasa wybrana kolumna target została ona przekształcona tak aby przechowywała tylko wartości 0 i 1 gdzie 1 będzie znaczyło chory .

2. Przetwarzanie i obróbka danych

Przetworzona została kolumna target która przechowywała wartości od 0 do 4, na taką która będzie miała wartości 0 i 1. Kolumna ta będzie służyć jako kolumna klasy.

Wykorzystana została biblioteka deducorrect i funkcja editWithRules. Poniżej kod odpowiedzialny za to.

```
setSickCorrection <- correctionRules(expression(  
  if(is.finite(target)){  
    if(target != 0)  
    {  
      if (target == 1 | target == 2 | target == 3 | target == 4) {  
        target <- 1  
      } else  
      {  
        target <- NA  
      }  
    }  
  }else  
  {  
    target <- NA  
  }  
}))
```

```
heart.disease.SickSet <- (correctWithRules(setSickCorrection,heart.disease))$corrected
```

Po edycji zostaje wyciągnięta tylko część kolumn z danymi (corrected) pomijając informacje o tym które zostały zmienione.

Funkcja zmieniała także każdą wartość która nie pasuje do wzorca na i w następnym kroku zostaną usunięte wiersze które zawierają NA w target, ponieważ interesują nas tylko te które mają wartość poprawną

```
heart.disease.Better <- subset(heart.disease.SickSet,is.finite(target))
```

Za pomocą biblioteki editrules zostało sprawdzone czy w bazie znajdują się jakiegokolwiek dane nie pasujące do wytycznych czy zakresów.

```
E <- editset(c(
  "age >0.0", "sex %in% c(0,1)", "cp %in% c(1,2,3,4)", "trestbps > 0",
  "chol > 0", "fbs %in% c(1,0)", "restecg %in% c(0,1,2,3)", "thalach > 0",
  "exang %in% c(0,1)", "oldpeak >=0", "slope %in% c(1,2,3)", "ca %in% c(0,1,2,3)",
  "thal %in% c(3,6,7)"
))
```

```
ve <- violatedEdits(E,heart.disease.Better)
```

Dzięki temu wiadomo dla których kolumn należy pisać funkcje uzupełniające dane bądź usuwające wiersze.

Po przeanalizowaniu kolumn wybrana została metoda kNN z biblioteki VIM aby na podstawie najbliższych sąsiadów uzupełnić luki w bazi. Kolumny te przechowują dane ściśle określone np.(1 lub 0), więc średnia nie była by oczekiwana w tym wypadku.

```
heart.knn <- kNN(heart.disease.Better)
```

3. Klasyfikatory i ich ewaluacja

Kolumną klasy jest target. Przechowuje ona wartość 0 i 1 czyli zdrowy i chory.

a)

Baza została podzielona na treningową i testową w stosunku 67/33.

Wylosowanie gdzie trafi który wiersz.

```
ind <- sample(2, nrow(heart.disease.ready), replace=TRUE, prob=c(0.67, 0.33))
```

i podzielenie bazy

```
heart.disease.train <- heart.disease.ready[ind==1, 1:14]
heart.disease.test <- heart.disease.ready[ind==2, 1:14]
```

b)

Użycie poszczególnych klasyfikatorów

Klasyfikator C4.5/ID3 znajduje się w bibliotece party

```
heart.disease.ctree <- ctree(target ~ age + sex + cp + trestbps + chol + fbs + restecg +
  thalach + exang + oldpeak + slope + ca + thal, data=heart.disease.train.ctree)
```

jego użycie polega na stworzeniu modelu drzewa powyższym poleceniem gdzie pierwszym argumentem jest wynik, który oczekujemy, elementy podane po ~ to już elementy wnioskowania, należy także podać bazę treningową.

Analiza danych następuje po przez następujące polecenia.

```
predicted.ctree <- predict(heart.disease.ctree, heart.disease.test[,1:13])  
trueSickValue <- heart.disease.test[,14]  
conf.matrix.ctree <- table(predicted.ctree,heart.disease.test[,14])  
conf.matrix.ctree  
accuracy.ctree <- sum(diag(conf.matrix.ctree))/sum(conf.matrix.ctree)  
accuracy.ctree
```

Polecenim predict uruchamia się model dla danych testowych.

Klasyfikator kNN

Do normalizacji wybrane zostały kolumny które posiadały pomiary i duże liczby aby nie powodować dużych różnic w dopasowaniach.

```
compute.model.norm.knn <- knn(heart.disease.norm.train[,1:13], heart.disease.norm.test[,1:13],  
cl=heart.disease.norm.train[,14], k = 3, prob=FALSE)
```

Tworzenie modelu który wyliczy wyniki na bazie testowej odbywa się po przez funkcję knn, w której podawana jest baza treningowa z kolumną target a testowa bez.

Analiza przebiega tak samo w C4.5/ID3

Klasyfikator Naive Bayes

```
modelHeart <- naiveBayes(target ~ ., data = heart.disease.train.NB)
```

Tworzenie modelu przebiega bardzo prosto podajemy kolumnę z docelowym wynikiem i nazwę bazy treningowej .

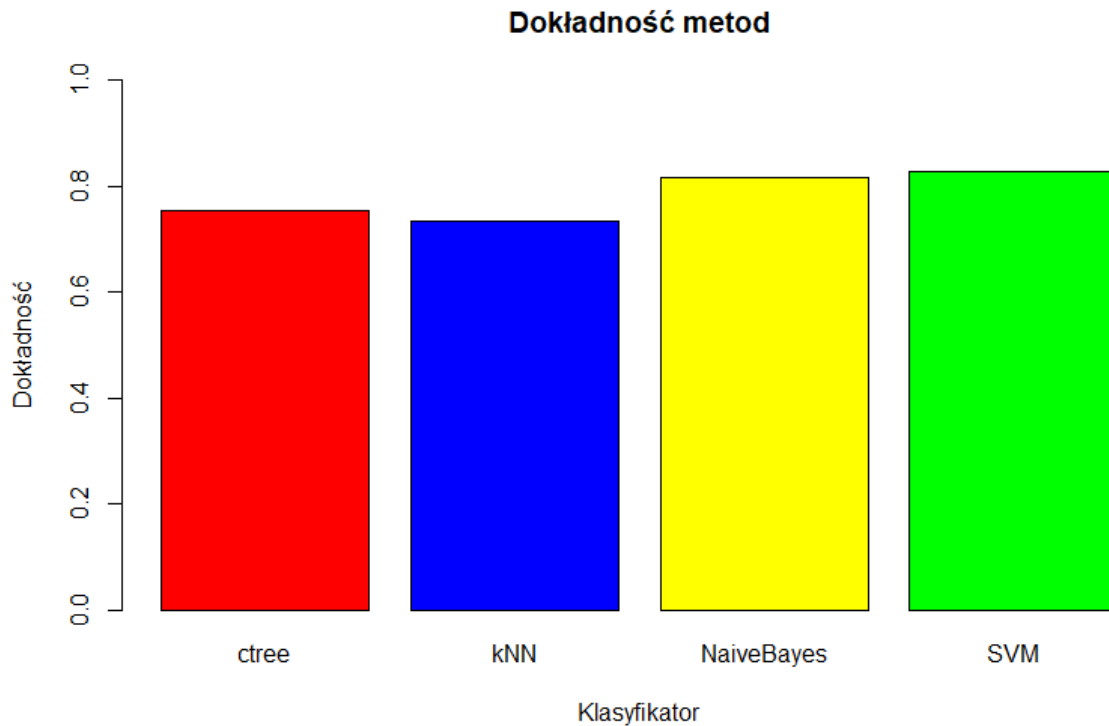
Analiza podobnie jak w pierwszym.

Dodatkowym klasyfikatorem jest SVM Support Vector Machines.

```
model.SVM <- svm(target ~ ., data = heart.disease.norm.train, type='C-classification',kernel='linear',  
scale=FALSE)
```

Porównanie wyników

Klasyfikator	Macierz błędu			Dokładność	TPR	FPR
Ctree	Real	Predicted		0.755102	0.875	0.4047619
		0	1			
		0 49 7				
kNN	Real	Predicted		0.7346939	0.7142857	0.2380952
		0	1			
		0 40 16				
naiveByes	Real	Predicted		0.8163265	0.8214286	0.1904762
		0	1			
		0 46 10				
SVM	Real	Predicted		0.8265306	0.8571429	0.2142857
		0	1			
		0 48 8				
	Real	Predicted				
		0	1			
		0 9 33				



c)

TP – oznacza liczbę osób, którym zgadł że są zdrowe

FP – oznacza liczbę pomylnych osób, które oznaczył jako chore a były zdrowe

FN – oznacza liczbę pomylnych osób, które oznaczył jako zdrowe a były chore

TN – oznacza liczbę osób, którym zgadł że są chore

TPR przedstawia odsetek poprawnie zdiagnozowanych(zdrowych) a TNR poprawnie zdiagnozowanych(chorych). Natomiast wartości FPR i FNR przedstawiają odsetki źle zdiagnozowanych. FNR jest dopełnieniem TPR. FPR dopełnieniem TNR

<p style="text-align: center;">Czułość, TPR</p> $\frac{\sum TP}{\sum TP + \sum FN}$	<p style="text-align: center;">FNR</p> $\frac{\sum FN}{\sum TP + \sum FN}$
<p style="text-align: center;">FPR</p> $\frac{\sum FP}{\sum FP + \sum TN}$	<p style="text-align: center;">Swoistość, SPC, TNR</p> $\frac{\sum TN}{\sum FP + \sum TN}$

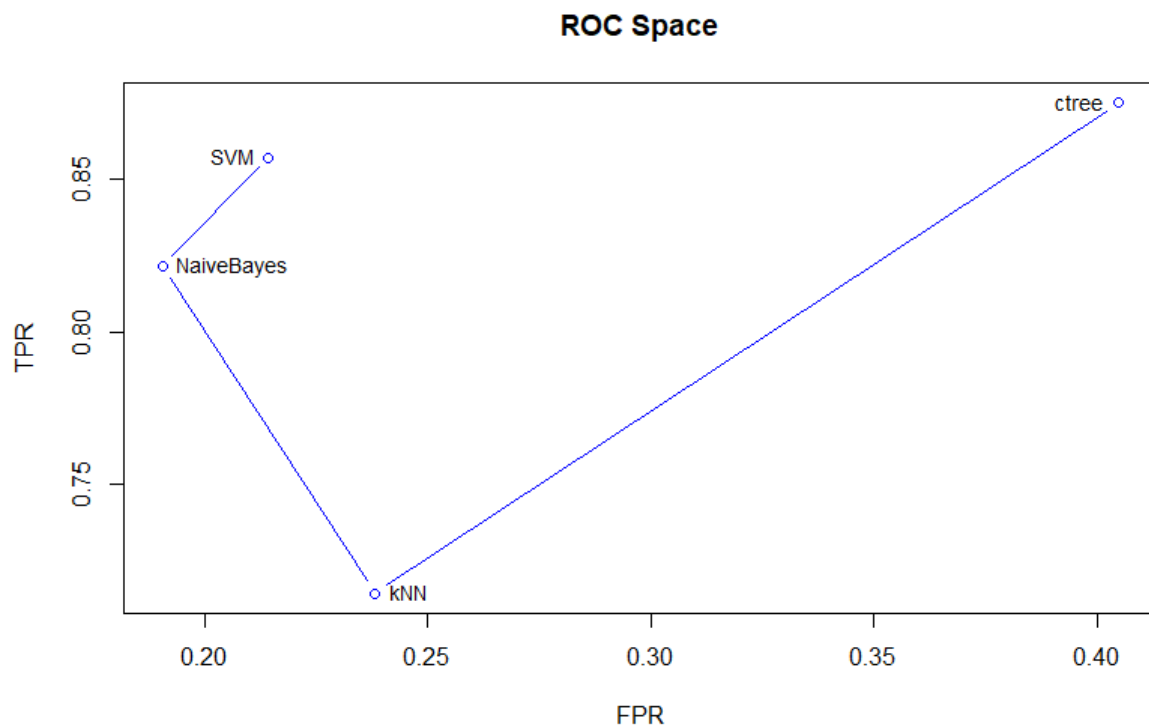
Błąd pierwszego stopnia u mnie to liczba osób źle sklasyfikowanych jako zdrowi a są chorzy.

Błąd drugiego stopnia u mnie to liczba osób zdrowych sklasyfikowanych jako chore.

Im więcej błędów drugiego stopnia to więcej osób mogło by być niepotrzebnie leczonych.

TNR i FPR przedstawiają skuteczność testu. A FNR i FPR nieskuteczność.

W mojej bazie i sposobie modelowania gorszy był by błąd pierwszego stopnia, ponieważ osoby chore nie były by leczone.



4. Grupowanie metodą k-średnich

Metoda k-średnich nie działa idealnie. Żadna z prób podzielenia na 2,3,4 czy 5 klastrów nie dała grup gdzie znajdziemy tylko chorych lub zdrowych. Przy podziale na dwie grupy skuteczność wyniosła ~ 0.835

Grupy		
	1	2
Real	0	13
	37	102

Czyli w jednej grupie 1 algorytm umieścić (zdrowi - wynioskowane) było 151 osób zdrowych i 37 chorych, a w grupie 2 (chorzy) 102 chorych i 13 zdrowych. Dokładność wydają się lepsza niż wśród Klasyfikatorów, ale trzeba wziąć pod uwagę że algorytm miał do dyspozycji całą bazę choć bez target.

Podsumowując można być zadowolonym wynikiem, ponieważ skuteczność jest porównywalna z najlepszy klasyfikatorem.

Dla przykładu wyniki dla 3,4 i 5 klastrow.

Ilość klastrow	Tabel					Skuteczność		
3	Real	Grupy				~0.8218		
		1	2	3				
		0	10	86	68			
		1	95	34	10			
4	Real	Grupy				~0.76567		
		1	2	3	4			
		0	8	23	66		67	
		1	76	20	7		36	
5	Real	Grupy					~0.76237	
		1	2	3	4	5		
		0	5	13	63	63		20
		1	56	28	27	7		21

Można dostrzec, że zwiększanie ilości klastrow pogarsza jakość grup, i bardzo złe grupy przeplatają się dobrymi grupami. Dla naszej bazy wydają się najkorzystniejszy podział na dwie grupy, lecz niektóre grupy z większej ilości klastrow mogą być korzystne.

Kod odpowiedzialny za grupowanie:

```
heart.disease.logPrep <- heart.disease.norm[,1:13]
heart.disease.log <- log(heart.disease.logPrep)
heart.disease.stand <- scale(heart.disease.logPrep, center=TRUE)
heart.disease.pca <- prcomp(heart.disease.stand)
heart.disease.final <- predict(heart.disease.pca)[,1:13]
```

Użycie algorytmu grupowania

```
k <- kmeans(heart.disease.final, 5)
```


5. . Reguły asocjacyjne

Użyta została biblioteka arules. Baza jest dość nieprzyjazna w tej metodzie poszukiwań ponieważ zawiera część danych pomiarowych, dla tego badania wykonywane są tylko na części kolumn.

Do badania wybrano kolumny : sex,cp,restecg,exang,thal,target a jako wynikowa wybrana sex,

Jedna asocjacja połączyła 4 informacje dla części osób

{restecg=0,exang=1,thal=3,target=1} => {sex=0} co znaczy że były kobiety u których powiązane były aspekty choroby . Co ciekawe oznacza to że osoby te zachorowały mając normalne wyniki poza zachorowaniem na anginę.

Kolejną ciekawą asocjacją w tym badaniu było {cp=1,restecg=0,thal=3,target=0} => {sex=0} gdzie przy normalnych wynikach i zachorowaniu na anginę nie wystąpiła choroba serca.

Do drugiego badania wybrano kolumny: cp, fbs, slope, target gdzie oczekiwaną jest target.

Ciekawym jest zestawienie cp=4 czyli braku bólu w klatce piersiowej i jednocześnie slope=4 czyli obniżenie częstoskurczu naczyniowego: {cp=4,slope=3} target=1

Porównie dwóch podobnych asocjacji

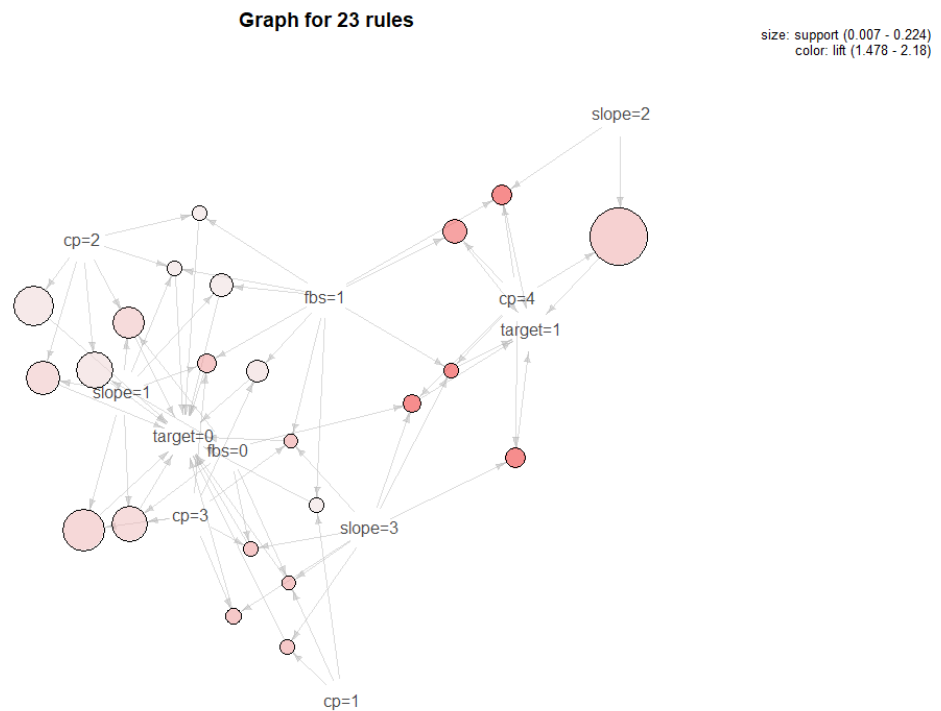
{cp=4,fbs=1,slope=2} => {target=1}

{cp=2,fbs=0,slope=1} => {target=0}

W pierwszym brak bólu jest powiązany z za dużym cukrem i płaskim częstoskurczem i daje nam chorobę.

W drugim mamy nietypową anginę i brak podwyższonego cukru oraz wzrosty częstoskurczu i brak choroby. Gdzie można by pomyśleć że wynik powinien być odwrotny.

Przykładowy graf asocjacji.



Kod odpowiedzialny za tworzenie reguł

```
rules <- apriori(mat1,  
  parameter = list(minlen=2, supp=0.005, conf=0.8),  
  appearance = list(rhs=c("sex=0", "sex=1"), default="lhs"),  
  control = list(verbose=F)  
)
```

Oraz przetworzenie ich.

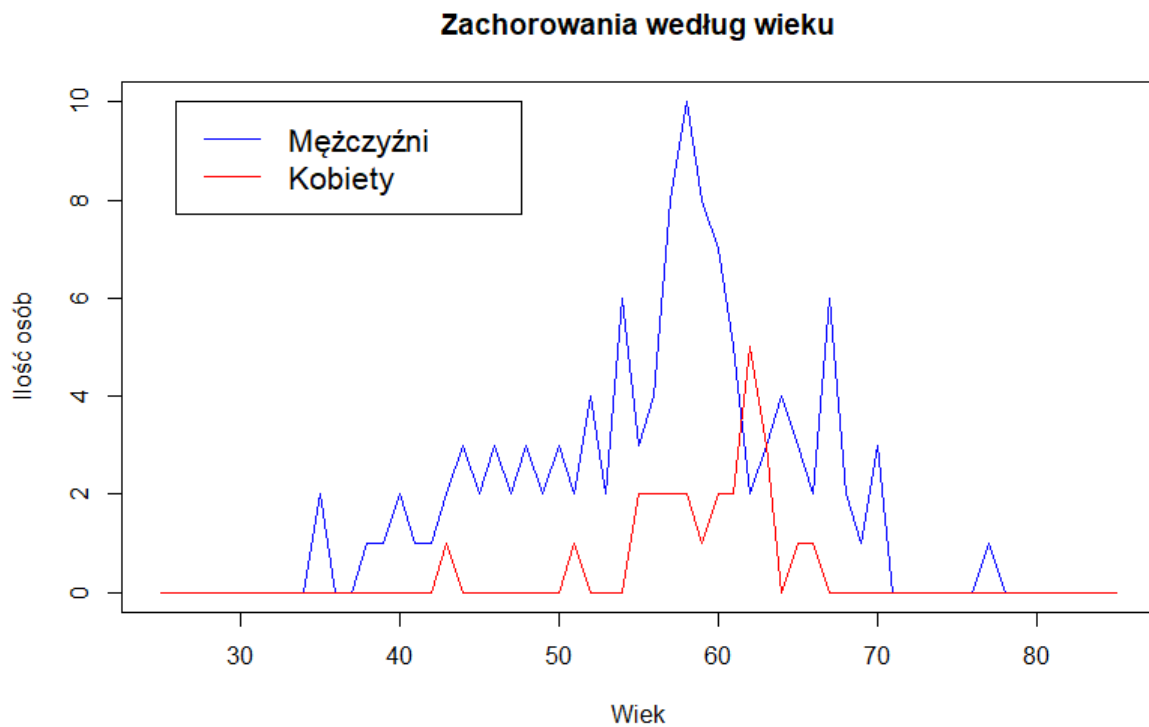
```
rules.sorted <- sort(rules, by="lift")  
subset.matrix <- is.subset(rules.sorted, rules.sorted)  
subset.matrix[lower.tri(subset.matrix, diag=T)] <- FALSE  
redundant <- colSums(subset.matrix, na.rm=T) >= 1  
rules.pruned <- rules.sorted[!redundant]
```

6. Dodatkowe wykresy

Liczba badanych mężczyzn 206 a kobiet 97.

Wiek minimalny kobiet 34, maksymalny 76.

Wiek minimalny mężczyzn 29, maksymalny 76

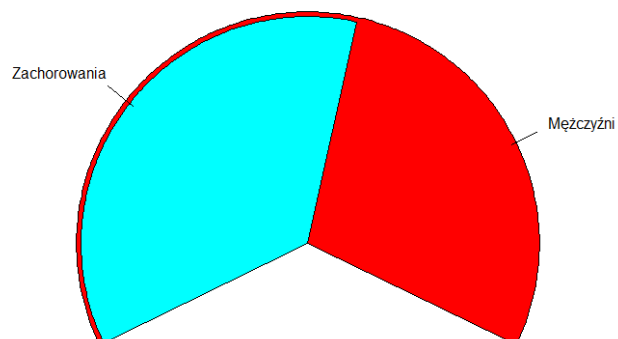


Wykres przedstawia zachorowania według wieku podzielone na płeć. Można zauważyć trend wzrostowy pomiędzy 50 i 60 rokiem życia dla kobiet i mężczyzn, który utrzymuje się do 70 roku życia. Celem wykresu było zauważenie trendów.

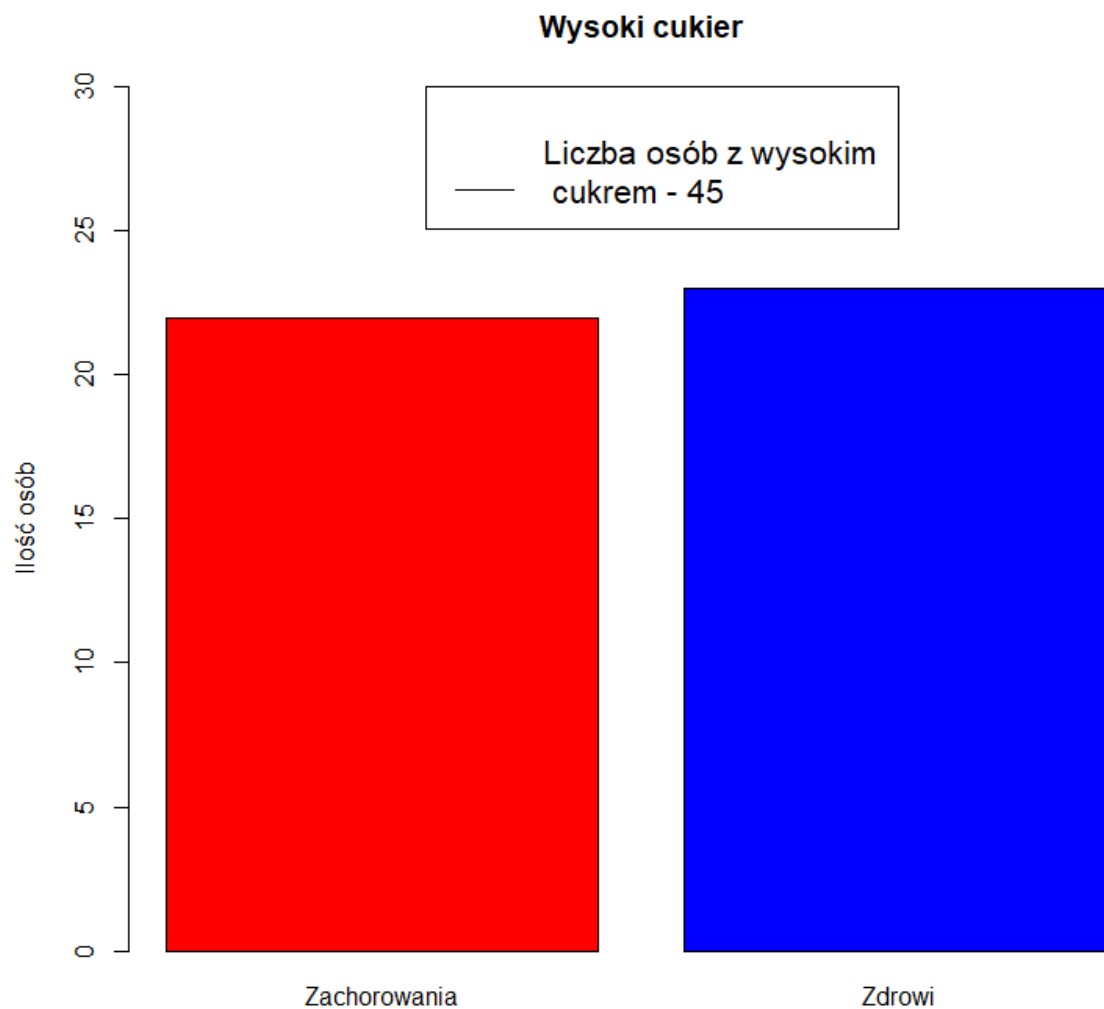
Procentowa zachorowalność kobiet



Procentowa zachorowalność mężczyzn



W grupie badanych gdy prześledzić procentowy udział zachorowań według płci, wychodzi że mężczyźni mają znacznie większą tendencję do zachorowania.



Powyższy wykres pokazuje stosunek zachorowań osób z wysokim cukrem $>120\text{mg/dl}$. Przewidywałem że ma on duży wpływ na choroby serca lecz w tej grupie badanych okazało się, połowa osób był chora a połowa nie. Więc to ma tylko częściowy wpływ. Rozkład 22 do 23 osób.

Dodatkowo przeprowadziłem badanie jak 3 klasyfikatory poradziły sobie bez normalizacji.

Klasyfikator	Macierz błędu (norm)	Macierz błędu (bez norm)	Dokładność
kNN	Real	Predicted	Z norm. 0.7346939 Bez norm. 0.6326531
		01	
		04016	
		11032	
naiveByes	Real	Predicted	Identyczna 0.8163265
		01	
		04610	
		1834	
SVM	Real	Predicted	Z norm. 0.8265306 Bez norm. 0.8061224
		01	
		0488	
		1933	

7. Podsumowanie

Klasyfikatory radzą sobie w miarę sprawnie biorąc nie tak wielką próbkę danych 303 rekordy. Głównym problem było przewidywanie choroby i na 80% można to przewidzieć w tym wypadku.

Najskuteczniejszy z klasyfikatorów był SVM, ale dobrze poradziła sobie też metoda k-średnich. Zaskakujące były asocjacje, które był niekiedy sprzeczne z tym co można by myśleć, i również to że ciężko na tej bazie operować asocjacjami. Najlepszą drogą było zmniejszenie zakresu większości kolumn do kilku wartości sztucznie nazwanych, tylko wkładając wyniki pomiarów do szuflad o określonym zakresie można by zakłamać nieco wyniki.