

# Analiza danych statystycznych dotyczących zachorowań na cukrzyce

**Adam Dyda**

*AGH, Wydział Informatyki Elektroniki i Telekomunikacji  
Rachunek prawdopodobieństwa i statystyka 2020/2021*

Kraków, 19 stycznia 2021

Adam Dyda

## Spis treści

<b>1</b>	<b>Opis danych</b>	<b>2</b>
<b>2</b>	<b>Czyszczenie danych</b>	<b>2</b>
<b>3</b>	<b>Analiza danych</b>	<b>4</b>
3.1	Analiza rozkładu danych . . . . .	4
3.2	Estymacja parametrów rozkładu . . . . .	5
3.3	Analiza wykresów pudełkowych . . . . .	9
<b>4</b>	<b>Analiza korelacji oraz zależności danych</b>	<b>12</b>
4.1	Analiza korelacji za pomocą współczynnika Pearsona . . . . .	12
4.2	Regresja liniowa oraz analiza korelacji . . . . .	13
<b>5</b>	<b>Wnioski</b>	<b>15</b>

## 1 Opis danych

Analizowane dane znajdują się na stronie <https://www.kaggle.com/uciml/pima-indians-diabetes-database>. Pochodzą one z instytutu *National Institute of Diabetes and Digestive and Kidney Diseases*. Dane dotyczą zachorowań na cukrzycę w grupie kobiet Indian Pima (*ang. Pima Indian*) powyżej 21 roku życia. Zawierają one informacje odnoszące się do kilku medycznych pomiarów i cech osób poddawanych analizie. Kolumny zbioru danych który będę analizował:

- Pregnancies - liczba ciąży
- Glucose - poziom glukozy we krwi mierzony podczas dwugodzinnego dostnego testu obciążenia glukoza
- BloodPressure - rozkurczowe ciśnienie krwi (mm Hg)
- SkinThickness - grubość fałdu skórniego na tricepsie (mm)
- Insulin - poziom insuliny we krwi
- BMI - wskaźnik masy ciała
- DiabetesPedigreeFunction (DPF) - funkcja określająca szanse na zachorowanie na podstawie historii przodków (ich zachorowań na cukrzycę) oraz relacji genetycznych z przodkami
- Age - wiek
- Outcome - cecha wskazująca czy dana osoba ma cukrzycę (0 lub 1)

## 2 Czyszczenie danych

Pierwszym etapem naszej analizy będzie sprawdzenie czy dane są kompletne i nie zawierają żadnych brakujących wartości, przyjrzymy się ich podsumowaniu.

Pregnancies		Glucose		BloodPressure		SkinThickness	
Min.	: 0.000	Min.	: 0.0	Min.	: 0.00	Min.	: 0.00
1st Qu.	: 1.000	1st Qu.	: 99.0	1st Qu.	: 62.00	1st Qu.	: 0.00
Median	: 3.000	Median	: 117.0	Median	: 72.00	Median	: 23.00
Mean	: 3.845	Mean	: 120.9	Mean	: 69.11	Mean	: 20.54
3rd Qu.	: 6.000	3rd Qu.	: 140.2	3rd Qu.	: 80.00	3rd Qu.	: 32.00
Max.	: 17.000	Max.	: 199.0	Max.	: 122.00	Max.	: 99.00
Insulin		BMI		DPF		Age	
Min.	: 0.0	Min.	: 0.00	Min.	: 0.0780	Min.	: 21.00
1st Qu.	: 0.0	1st Qu.	: 27.30	1st Qu.	: 0.2437	1st Qu.	: 24.00
Median	: 30.5	Median	: 32.00	Median	: 0.3725	Median	: 29.00
Mean	: 79.8	Mean	: 31.99	Mean	: 0.4719	Mean	: 33.24
3rd Qu.	: 127.2	3rd Qu.	: 36.60	3rd Qu.	: 0.6262	3rd Qu.	: 41.00
Max.	: 846.0	Max.	: 67.10	Max.	: 2.4200	Max.	: 81.00

```

Outcome
Min.    :0.000
1st Qu.:0.000
Median  :0.000
Mean    :0.349
3rd Qu.:1.000
Max.    :1.000

```

Jak widac dane nie zawieraja brakujacych wartosci (mozna to zaobserwować poprzez brak wpisu "NA's" w ktorejkoľwiek z kolumn). Mozemy jednak zauwazyc ze niektore z kolumn posiadaja 0 jako wartosc minimalna. Zastanowmy sie teraz czy dla wszystkich kolumn wartosc 0 moze byc rzeczywiscie zmierzona i jest mozliwa do uzyskania. Poziom glukozy, cislienie krwi, poziom insuliny, grubosc skory oraz BMI nie moga przyjmowac wartosci 0. Sprawdzmy teraz ile jest wartosci 0 w calym zbiorze danych.

```
> colSums(df == 0)
```

Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin
111	5	35	227	374
BMI	DPF	Age	Outcome	
11	0	0	500	

Wywnioskowac z tego mozemy ze w zbiorze brak danych jest oznaczany wartoscia 0 (taki wniosek przyjmujemy dla kolumn wzczesniej wymienionych dla ktorych ta wartosc "nie ma sensu" w pozostalych kolumnach uznajemy wartosc 0 jako wartosc poprawna)

Brakujace dane nalezy odrzucic lub uzupelnic, odrzucenie wierszy z brakujacymi danymi w tym wypadku wiaze sie z utrata bardzo duzej ilosci informacji poniewaz dla kolumny *Insulin* sa to az 374 wartosci, rozwazyc mozna takze odrzucenie calej kolumny jednak informacja o poziomie insuliny jest dla nas zbyt istotna przy badaniu zachorowan na cukrzyce.

Zdecydowalem sie na uzupelnienie brakujacych danych odpowiednimi wartosciami w zaleznosci od kolumny. W przypadku kolumn z dostatecznie mala skosnoscia rozkladu brakujace dane bede uzupelnial srednia arytmetyczna z pozostalych poprawnych danych w tej kolumnie, w przypadku rozkladu z duza skosnoscia uzyje mediany do uzupelnienia danych. Sprawdzmy skosnosc wszystkich kolumn w zbiorze:

```
> skewness(df)
```

Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin
0.8999119	0.1734140	-1.8400052	0.1091588	2.2678105
BMI	DPF	Age	Outcome	
-0.4281433	1.9161592	1.1273893	0.6337757	

Za skosnosc dostateczna do uzupełnienia danych srednia arytmetyczna przyjmuje wartosc wieksza od -0.5 oraz mniejsza od 0.5. Oczywiscie uzupełnienie brakujacych danych srednia lub mediana wprowadza przeklamania z ktorych nalezy zdawac sobie sprawe, dlatego tez w dalszej czesci analizy jezeli bedzie to potrzebne beda uzywal danych oryginalnych z odrzuconymi danymi brakujacymi (dotyczy to glownie kolumny *Insulin*). Taki zabieg zostanie odpowiednio zasygnalizowany wzczesniej.

### 3 Analiza danych

Zajmiemy sie teraz analiza rozkladu danych w naszym zbiorze, przyjrzymy sie histogramom oraz dokonamy estymacji parametrow naszych rozkladow za pomoca stymatora MLE (*ang. Maximum likelihood estimator*). Nastepnie przejdziemy do analizy boxplotow oraz powiemy o wartosciach odstajacych w naszych danych.

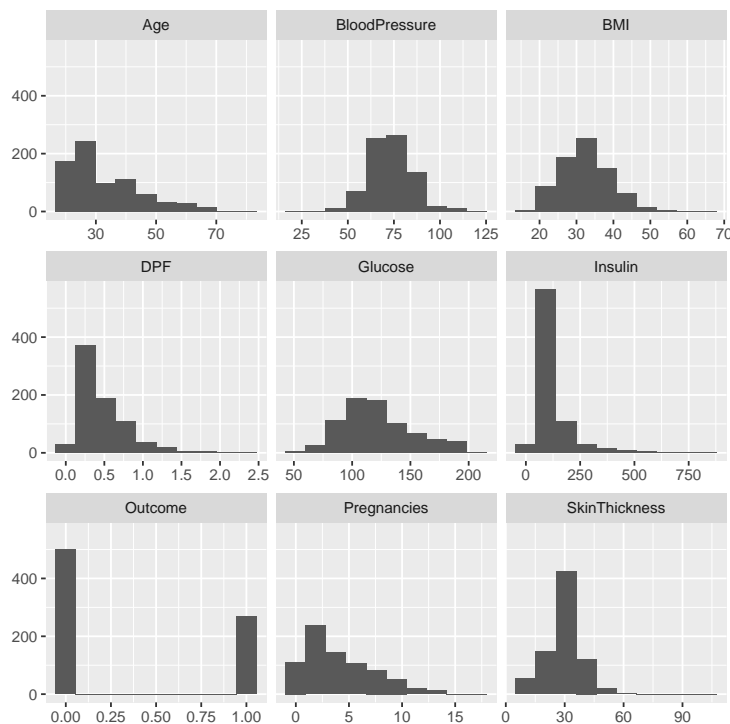
#### 3.1 Analiza rozkladu danych

Analize zaczniemy od ponownego przeanalizowania rezultatu polecenia *summary* tym razem juz dla uzupełnionych danych:

<b>Pregnancies</b>	<b>Glucose</b>	<b>BloodPressure</b>	<b>SkinThickness</b>
Min. : 0.000	Min. : 44.00	Min. : 24.00	Min. : 7.00
1st Qu.: 1.000	1st Qu.: 99.75	1st Qu.: 64.00	1st Qu.:25.00
Median : 3.000	Median :117.00	Median : 72.20	Median :29.00
Mean : 3.845	Mean :121.66	Mean : 72.41	Mean :29.11
3rd Qu.: 6.000	3rd Qu.:140.25	3rd Qu.: 80.00	3rd Qu.:32.00
Max. :17.000	Max. :199.00	Max. :122.00	Max. :99.00
<b>Insulin</b>	<b>BMI</b>	<b>DPF</b>	<b>Age</b>
Min. : 14.0	Min. :18.20	Min. :0.0780	Min. :21.00
1st Qu.:121.5	1st Qu.:27.50	1st Qu.:0.2437	1st Qu.:24.00
Median :125.0	Median :32.30	Median :0.3725	Median :29.00
Mean :140.7	Mean :32.46	Mean :0.4719	Mean :33.24
3rd Qu.:127.2	3rd Qu.:36.60	3rd Qu.:0.6262	3rd Qu.:41.00
Max. :846.0	Max. :67.10	Max. :2.4200	Max. :81.00
<b>Outcome</b>			
Min. :0.000			
1st Qu.:0.000			
Median :0.000			
Mean :0.349			
3rd Qu.:1.000			
Max. :1.000			

Na podstawie tych informacji mozna zaobserwowac ze w przypadku kazdej z cech rozklad jest prawostronnie skosny (srednia arytmetyczna jest wieksza od mediany). W przypadku kolumn *Insulin*, *Pregnancies*, *DPF*, *SkinThickness* widzimy duza roznice pomiedzy trzecim kwartylem a wartoscia maksymalna co

moze oznaczac jedna lub potencjalnie wiele wartosci odstajacych. Przyjrzyjmy sie teraz histogramom naszych danych:



Na powyzzszych histogramach widzimy ze kolumny *BloodPressure*, *Glucose* oraz *BMI* posiadaja rozklad ktory przypomina rozklad normalny. A wiec wlasnie tymi cechami zajmiemy sie w kolejnej sekcji podczas estymacji parametrow rozkladu.

Warto takze zwrocic uwage na rozklad danych w kolumnie *Outcome* i zauwazyc ze nie sa one rownomiennie rozlozone. Sprawdzmy dokladna liczbe osob chorych i zdrowych.

```
> table(df$Outcome)
```

```
0    1
500 268
```

Jak widac liczba osob zdrowych jest prawie dwa razy wieksza od liczby osob chorych, co moze miec znaczenie w dalszej analizie.

### 3.2 Estymacja parametrow rozkladu

Teraz zajmiemy sie estymacja parametrow  $\mu$  oraz  $\sigma^2$  rozkladow trzech cech ktore wybralismy w poprzednim punkcie tj. *BloodPressure*, *Glucose* oraz *BMI*. Estymacji bedziemy dokonywac za pomoca estymatora największej wiarygodności

MLE. Na początku zdefiniujemy sobie funkcje wiarygodności, w naszym wypadku jest to *Negative Likelihood function*, co oznacza że będziemy minimalizować te funkcje.

```
> NLLGlucose <- function(theta0,theta1) {
+   -sum ( -0.5* log(theta1*2*pi) - 0.5*( df$Glucose- theta0)^2/theta1 )
+ }
```

Analogicznie definiujemy funkcje dla kolumn *BloodPressure* oraz *BMI*. Następnie minimalizujemy funkcje wiarygodności aby otrzymać odpowiednie parametry, używamy do tego funkcji *mle* z pakietu *stats4*

*theta0* oznacza parametr  $\mu$  natomiast *theta1* oznacza  $\sigma^2$

```
> Glucoseest <- stats4::mle(minuslog=NLLGlucose, start=list(theta0=100,theta1=900))
> coef(Glucoseest)
```

```
theta0    theta1
121.6624  913.1354
```

```
> BMIest <- stats4::mle(minuslog=NLBMI, start=list(theta0=30,theta1=60))
> coef(BMIest)
```

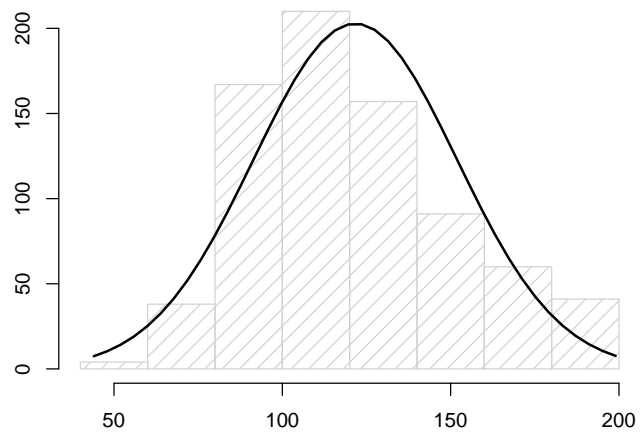
```
theta0    theta1
32.45521  47.20657
```

```
> BPest <- stats4::mle(minuslog=NLBP, start=list(theta0=70,theta1=375))
> coef(BPest)
```

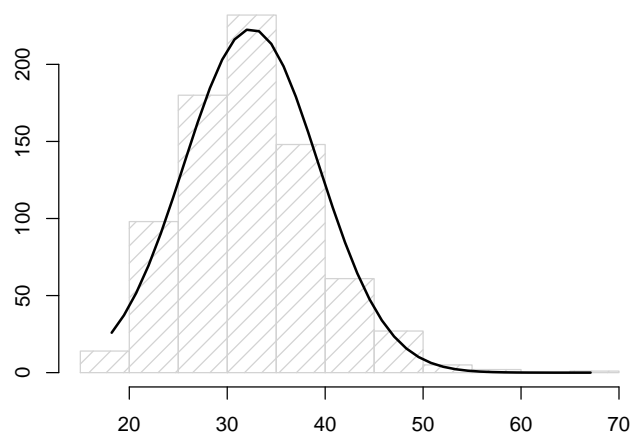
```
theta0    theta1
72.40497 145.92579
```

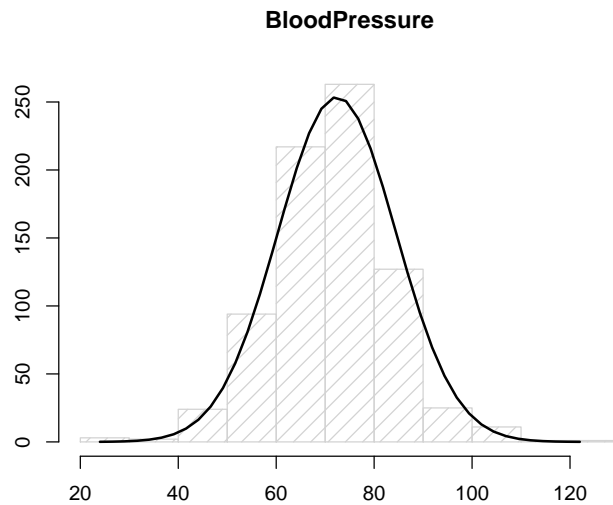
Na podstawie wyestymowanych informacji narysujmy wykresy rozkładów normalnych z otrzymanymi parametrami na histogramach każdej z cech.

**Glucose histogram**



**BMI histogram**





Wystymowane rozkłady wyglądają na bliskie prawdy tj. takie których parametry bliskie są rzeczywistym parametrom rozkładu. Porównajmy więc wartości rzeczywiste z wystymowanymi.

Rzeczywiste wartości parametrów dla kolumny *Glucose*

```
> mean(df$Glucose)
```

```
[1] 121.6562
```

```
> var(df$Glucose)
```

```
[1] 926.4892
```

Estymowane wartości parametrów dla kolumny *Glucose*  
 $\theta_0$  oznacza parametr  $\mu$  natomiast  $\theta_1$  oznacza  $\sigma^2$

```
theta0  theta1
121.6624 913.1354
```

Rzeczywiste wartości parametrów dla kolumny *BMI*

```
> mean(df$BMI)
```

```
[1] 32.45521
```

```
> var(df$BMI)
```

```
[1] 47.26806
```



Estymowane wartosci parametrow dla kolumny *BMI*  
*theta0* oznacza parametr  $\mu$  natomiast *theta1* oznacza  $\sigma^2$

```
theta0    theta1  
32.45521  47.20657
```

Rzeczywiste wartosci parametrow dla kolumny *BloodPressure*

```
> mean(df$BloodPressure)
```

```
[1] 72.40518
```

```
> var(df$BloodPressure)
```

```
[1] 146.3216
```

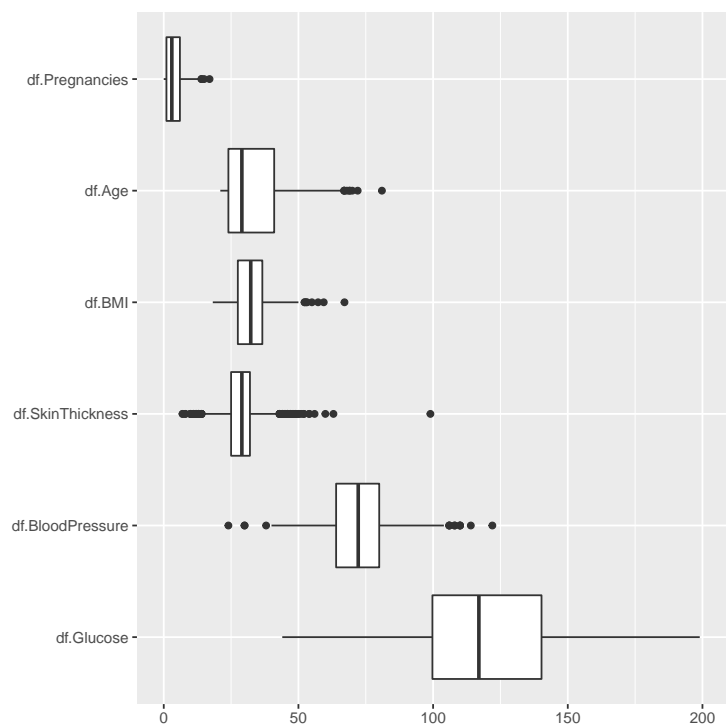
Estymowane wartosci parametrow dla kolumny *BloodPressure*  
*theta0* oznacza parametr  $\mu$  natomiast *theta1* oznacza  $\sigma^2$

```
theta0    theta1  
72.40497  145.92579
```

Estymowane wartosci sa bardzo bliskie prawdy, mozna wiec stwierdzic ze estymacja za pomoca metody najmniejszych kwadratow okazala sie wlasciwa w tym przypadku.

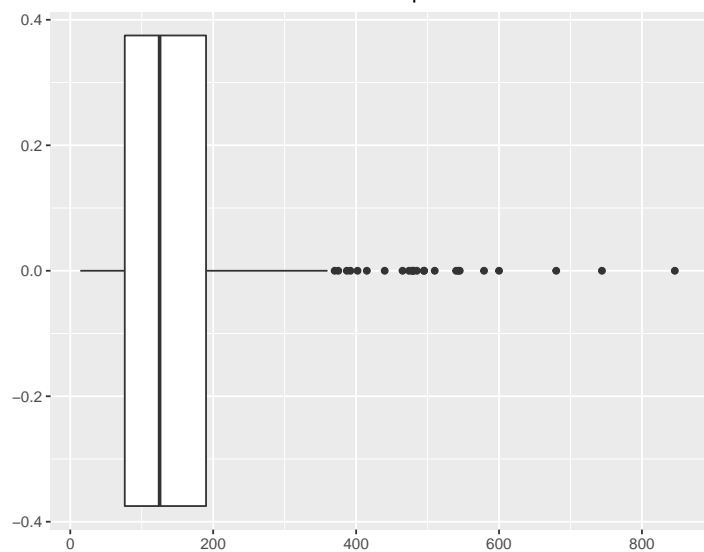
### 3.3 Analiza wykresow pudelkowych

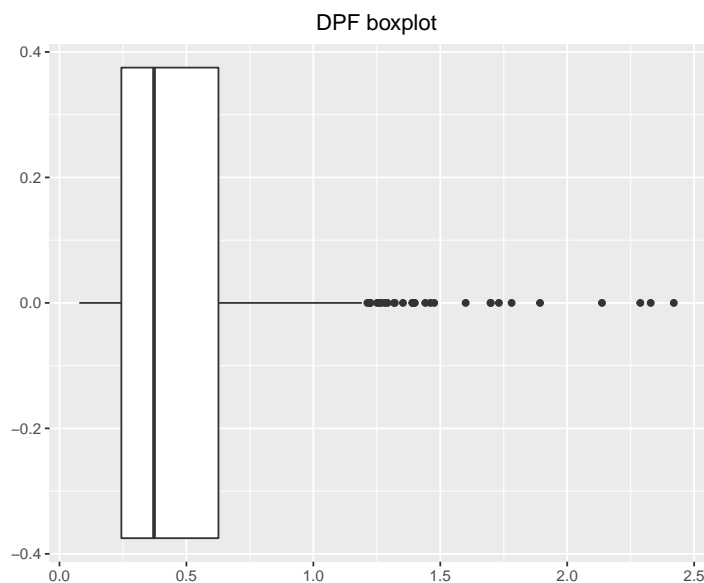
Dalsza czescia eksploracji naszych danych bedzie narysowanie oraz analiza wykresow pudelkowych kazdej z cech. Wykresy rysuje w trzech oddzielnych grupach poniewaz wartosci sa z roznych rzadow wielkosci i narysowane na jednym wykresie bylyby nieczytelne.



Podczas rysowania wykresu pudełkowego kolumny *Insulin* zostały użyte dane z usuniętymi wierszami których wartość w oryginalnym zbiorze wynosiła 0.

Insulin boxplot





Na powyższych wykresach możemy zauważyć znaczącą ilość wartości odstających w kolumnach *Insulin*, *DPF* oraz *SkinThickness*, co może mieć wpływ na zakłamanie podczas używania regresji liniowej w dalszej części analizy. Wynika to z tego, że model regresyjny jest bardzo czuły na wartości odstające.

Sprawdźmy dokładną liczbę wartości odstających w tych kolumnach. Za wartości odstające przyjmujemy wartości mniejsze od  $Q1 - 1.5IQR$  oraz większe od  $Q3 + 1.5IQR$  gdzie  $Q1, Q3$  to odpowiednio wartości pierwszego i trzeciego kwartyła, a  $IQR$  to rozstęp cwiartkowy, czyli różnica między trzecim i pierwszym kwartyłem.

```
> length(boxplot(cleandf$Insulin)$out)
[1] 24

> length(boxplot(df$DPF)$out)
[1] 29

> length(boxplot(df$SkinThickness)$out)
[1] 87
```

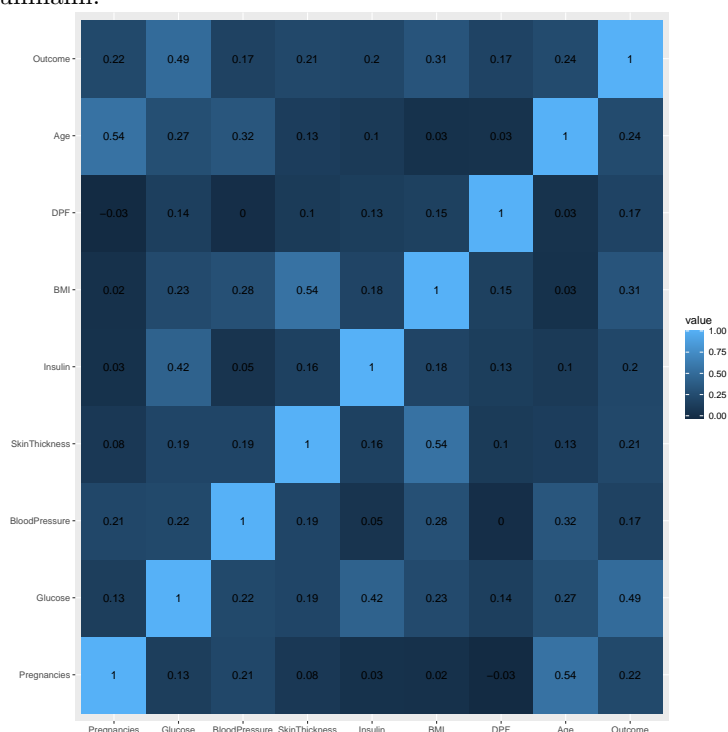
Warto zauważyć, że wartość dla kolumny *SkinThickness* jest zawyżona, ponieważ używamy tutaj danych z uzupełnionymi brakującymi wartościami.

## 4 Analiza korelacji oraz zaleznosci danych

W tym kroku zajmiemy sie analiza zaleznosci pomiedzy danymi za pomoca wspolczynnika korelacji Pearsona. W szczegolnosci zajmiemy sie badaniem korelacji pomiedzy kolumna *Outcome* a pozostalymi cechami w celu ustalenia ktore z cech maja najwieksza korelacje z zachorowaniem na cukrzyce.

### 4.1 Analiza korelacji za pomoca wspolczynnika Pearsona

W celu zwizualizowania korelacji miedzy kolumnami wykorzystam wykres typu heatmap ktory w jednym miejscu przedstawia zaleznosci pomiedzy wszystkimi kolumnami.

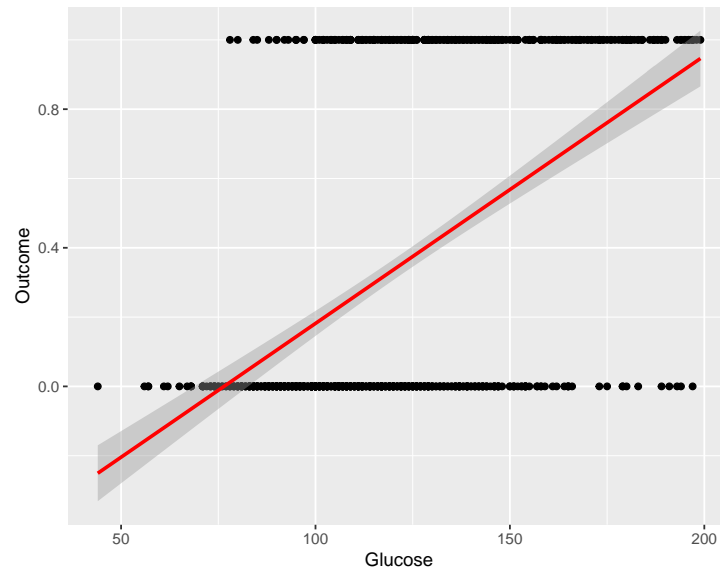


Na pierwszy rzut oka mozemy zobaczyc np. wysoka korelacje wieku z liczba ciąż lub wieku z cislaniem krwi. Sa to oczywiste zaleznosci ktorymi nie bedziemy sie glebiej zajmowac jednakze warto na nie zwrocic uwage. Wartosci ktore nas interesuja jest to korelacja zachorowan na cukrzyce z reszta kolumn, w tym wypadku dwie najwieksze wartosci czyli korelacje z poziomem glukozy we krwi oraz wskaźnikiem BMI, interesuje nas takze korelacja zachorowan na cukrzyce z poziomem insuliny we krwi poniewaz jest to bardzo czesto badana cecha wsrod chorych na te chorobe.

## 4.2 Regresja liniowa oraz analiza korelacji

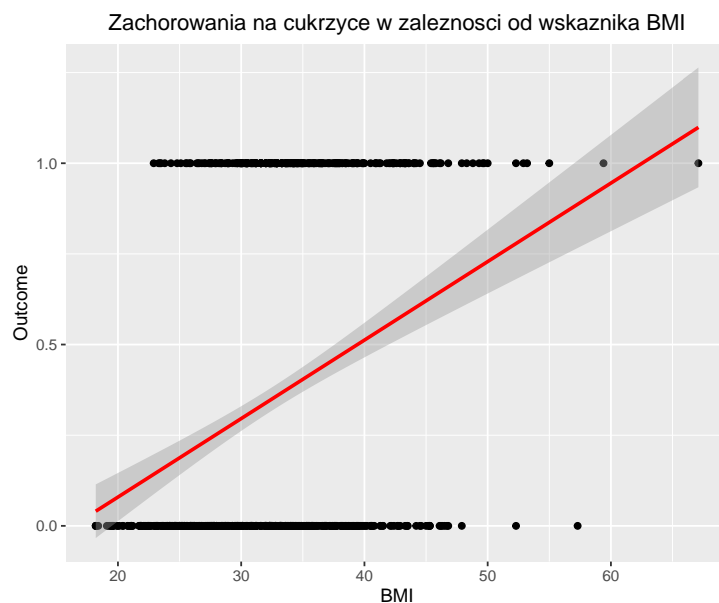
W tym punkcie użyjemy modelu regresji liniowej do wizualizacji zależności pomiędzy wybranymi cechami. Pierwszym wykresem będzie wykres zależności poziomu glukozy do zachorowań na cukrzyce do którego dopasowujemy prostą za pomocą regresji liniowej.

Zachorowania na cukrzyce w zależności od poziomu glukozy we krwi

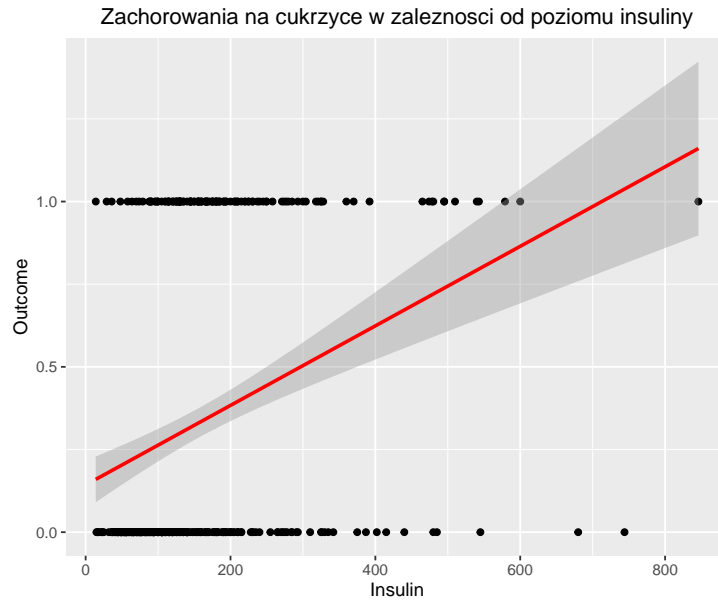


Jak widzimy korelacja jest pozytywna co oznacza że im większy poziom glukozy we krwi tym większa szansa na to że osoba jest chora na cukrzyce. Jest to wynik którego się spodziewaliśmy ponieważ często wysoki poziom cukru jest właśnie indykatorem cukrzycy typu 1 lub 2.

Następnie narysujmy za pomocą modelu regresji liniowej wykres zależności wskaźnika BMI do zachorowań na cukrzyce



W tym wypadku korelacja także jest pozytywna co oznacza że im większy wskaźnik BMI tym większa szansa zachorowania na cukrzyce. Jest to wniosek który potwierdza ogólne przekonanie że osoby otyłe częściej chorują na cukrzyce. Wykresem kolejnej zależności będzie zależność poziomu insuliny we krwi do zachorowań na cukrzyce. W tym wypadku używam zbioru danych z usuniętymi wierszami których wartość poziomu insuliny w oryginalnym zbiorze wynosiła 0.



Zależność w tym wypadku także jest pozytywna co jest zaskakującą obserwacją, ponieważ spodziewamy się, że osoby chore na cukrzyce mają niski poziom insuliny. W naszym wypadku można wyciągnąć dalsze wnioski z tej obserwacji, biorąc pod uwagę, że osoby chore na cukrzyce typu 1 mają niski poziom insuliny, natomiast osoby chore na cukrzyce typu 2 mogą mieć ten poziom insuliny podwyższony (co nie jest w tym wypadku regułą) m.in. dlatego że hiperinsulinemia może prowadzić do rozwoju cukrzycy typu 2. Możemy wywnioskować z danych osób chorych w naszym zbiorze pochodzą od osób w większości (lub wszystkich) chorych na cukrzyce typu 2.

Warto zaznaczyć, że powyższe dopasowania linii za pomocą regresji liniowej są niedokładne z kilku powodów. Regresja liniowa sprawdza się najlepiej w wypadku, w którym cechy pochodzą z rozkładu normalnego, w naszym zbiorze kolumna *Outcome* przyjmuje tylko dwie wartości, więc zanika to poziom skuteczności naszego modelu. Kolejnym powodem jest już wcześniej wspomniany problem z wartościami odstającymi, na które model regresji liniowej jest bardzo czuły, jednym z rozwiązań tego problemu mogłoby być użycie regularyzowanej regresji liniowej, czyli regresji z użyciem parametru penalizującego. (*Ridge Regression*, *Lasso Regression*)

## 5 Wnioski

Analiza danych zakończyła się sukcesem, udało nam się sprawnie wyczyścić dane i doprowadzić je do stanu, w którym są użyteczne i podatne do analizowania. Dalsza analiza pozwoliła nam odkryć asymetryczności i nierównomierności w naszych danych. Kolejny etap, czyli estymacja parametrów na pomocą metody

MLE także dała zamierzony efekt, ponieważ udało nam się sprawnie i zadowalająco dokładnie wyestymować parametry badanych rozkładów. Przyjrzelismy się wykresom pudełkowym co pozwoliło określić które z cech posiadają wartości odstające które są potencjalnym problemem w dalszych etapach analizy. Następnie przeszliśmy do analizy korelacji między cechami i wyciągania wniosków na temat zależności między nimi.

Główną zależnością którą chcieliśmy zbadać jest korelacja zachorowań na cukrzycę z cechami osoby potencjalnie chorej. Analizując nasze wyniki można dojść do wniosku że badanie zakończyło się sukcesem, ponieważ zauważyliśmy kilka pozytywnych korelacji z zachorowaniami na cukrzycę m.in. poziom glukozy we krwi oraz wskaźnik BMI co potwierdza powszechną opinię że poziom glukozy należy kontrolować oraz że osoby otyłe częściej chorują na cukrzycę. Wyjątkowo ciekawy jest wniosek płynący z analizy korelacji insuliny z zachorowaniem na cukrzycę. Z tej zależności wywnioskowaliśmy że osoby chore w naszym zbiorze danych chorują głównie na cukrzycę typu 2, za tą tezę przemawiać może także fakt że dane są pobrane od zamkniętej grupy osób o określonej charakterystyce co predysponuje taką grupę do chorowania na węższy zakres chorób.

Po zakończeniu analizy danych należy zdawać sobie sprawę z niedokładności i przekłaman które w niej występują oraz które wprowadziliśmy poprzez uzupełnienie brakujących danych. Pod uwagę należy wziąć także to że zarówno współczynnik korelacji Pearsona jak i model regresji liniowej swoją najlepszą efektywność uzyskują dla danych ciągłych o rozkładzie normalnym, jednak nie wszystkie nasze cechy takie właściwości posiadały. W szczególności kolumna *Outcome* jest cechą jakościową (tj. niemierzalna, przyjmująca wartości 0 lub 1). Przekłamania wprowadza także liczba wartości odstających w niektórych kolumnach. Warto się zastanowić czy w analizie tych danych nie lepszym rozwiązaniem byłoby użycie innych modeli statystycznych, mogłoby to pozwolić nam w uzyskaniu dokładniejszych wyników i uniknąć przekłaman i niedokładności w niektórych aspektach analizy.