

```
set.seed(42)
```

Arbeitsblatt 9

Aufgabe 1

Unterschiede hinsichtlich

- der Voraussetzungen an die Zielvariable: Bei Diskriminanzanalyse und KNN auch Multilabel-Classification möglich, bei logistischer Regression nur Binary-Classification.
- den Voraussetzungen an die Einflussvariablen: Bei logistischer Regression und KNN keine Normalverteilungsannahme der Einflussvariablen nötig, bei Diskriminanzanalyse schon.
- der resultierenden Klassifikationsgrenzen: Bei logistischer Regression und KNN sind Klassifikationsgrenzen interpretierbar, da es sich um geschätzte Wahrscheinlichkeiten handelt, bei Diskriminanzanalyse ist eine Interpretation nicht möglich.
- den Konsequenzen einer Standardisierung von Einflussvariablen: Bei KNN ist eine Standardisierung nötig (nur hilfreich?), da Distanzmaße skalenabhängig sind, logistische Regression und Diskriminanzanalyse sind skalenunabhängig.

Aufgabe 2

a)

```
load("diab.RData")

log_reg <- glm(diabetes~pgc+bmi, family=binomial(link='logit'), data=diab)

pred = predict.glm(log_reg, data.frame(pgc=120, bmi=27), type="response")[1]
pred = as.vector(pred)

logit = function(L) {
  return(1/2*(1+tanh(L/2)))
}

pred2 = as.vector(logit(log_reg$coefficients[1]+log_reg$coefficients[2]
                        *120+log_reg$coefficients[3]*27))
delta = round((pred - pred2)^2, 4)
delta
#> [1] 0

pred
#> [1] 0.2254225
pred2
#> [1] 0.2254225
```

b)

```
summary(log_reg)
#>
#> Call:
#> glm(formula = diabetes ~ pgc + bmi, family = binomial(link = "logit"),
#>      data = diab)
```

```
#>
#> Deviance Residuals:
#>      Min       1Q   Median       3Q      Max
#> -2.1803  -0.7752  -0.4693   0.7758   3.0164
#>
#> Coefficients:
#>              Estimate Std. Error z value Pr(>|z|)
#> (Intercept) -7.515639   0.605236 -12.418  < 2e-16 ***
#> pgc          0.035169   0.003289  10.694  < 2e-16 ***
#> bmi          0.076334   0.013338   5.723 1.05e-08 ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> (Dispersion parameter for binomial family taken to be 1)
#>
#>      Null deviance: 993.48  on 767  degrees of freedom
#> Residual deviance: 771.40  on 765  degrees of freedom
#> AIC: 777.4
#>
#> Number of Fisher Scoring iterations: 4
```

Die Koeffizientenschätzung der Variablen `pgc` ist $\beta = 0.035169$, ein positiver Wert. Dies bedeutet, dass ein Anstieg der Glukosekonzentration mit einer Erhöhung der Wahrscheinlichkeit verbunden ist, an Diabetes erkrankt zu sein.

c)

Ein wichtiges zu verstehendes Konzept für die Interpretation der logistischen Beta-Koeffizienten ist die Odds Ratio. Die Odds Ratio misst die Zuordnung zwischen einer Prädiktorvariablen (x) und der Ergebnisvariablen (y). Es stellt das Verhältnis der Chancen dar, dass ein Ereignis eintritt (Ereignis = 1) bei Vorhandensein des Prädiktors x ($x = 1$), verglichen mit den Chancen des Ereignisses, das in Abwesenheit dieses Prädiktors eintritt ($x = 0$).

Für einen gegebenen Prädiktor (z.B. x_1) entspricht der zugehörige Beta-Koeffizient (b_1) in der logistischen Regressionsfunktion dem Logarithmus der Odds Ratios für diesen Prädiktor.

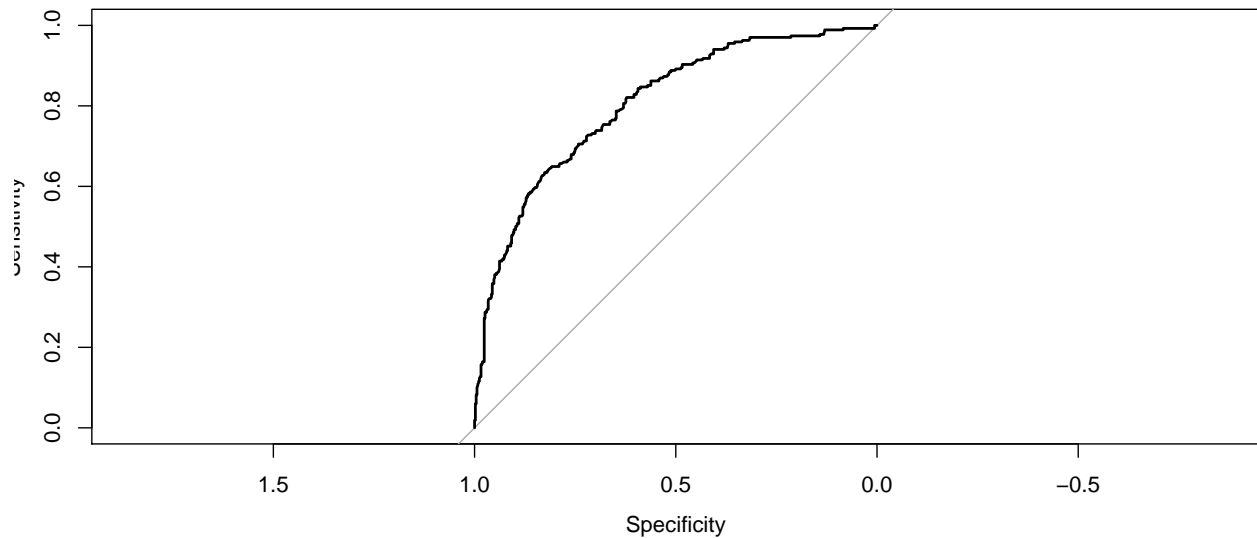
Wenn die Odds Ratio 2 ist, dann sind die Chancen, dass das Ereignis eintritt (Ereignis = 1), doppelt so hoch, wenn der Prädiktor x vorhanden ist ($x = 1$) als für ein Fehlen des Prädiktors ($x = 0$).

So beträgt beispielsweise der Regressionskoeffizient für BMI 0.0763342. Dies deutet darauf hin, dass eine Erhöhung des BMI um 5 Einheiten die Wahrscheinlichkeit, an Diabetes erkrankt um 0.32% erhöht $1 - 1/\exp(5\beta_{\text{BMI}})$.

d)

```
packageTest("pROC")

plot.roc(diab$diabetes, diab$pgc+diab$bmi)
```



e)

```

test <- NULL
test$y.true <- diab$diabetes
test$y.pred <- predict.glm(log_reg, type="response")

test$y.pred[test$y.pred > 0.5] = 1
test$y.pred[test$y.pred <= 0.5] = 0

confusion_matrix <- table(test)
confusion_matrix
#>      y.pred
#> y.true  0   1
#>    0 445  55
#>    1 126 142

misclass_rate <- 1-sum(diag(confusion_matrix))/sum(confusion_matrix)

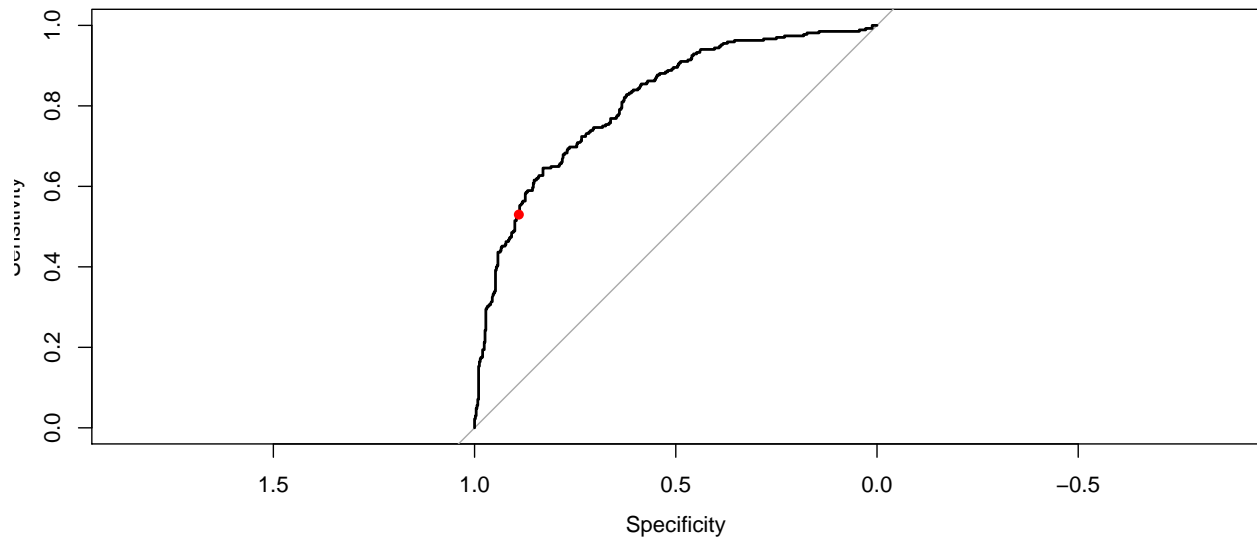
sensitivity <- confusion_matrix[2,2]/(confusion_matrix[2,1] + confusion_matrix[2,2])

specificity <- confusion_matrix[1,1]/(confusion_matrix[1,1] + confusion_matrix[1,2])

misclass_rate
#> [1] 0.2356771
sensitivity
#> [1] 0.5298507
specificity
#> [1] 0.89

plot.roc(diab$diabetes, predict.glm(log_reg, type="response"))
points(specificity, sensitivity, col="red", pch = 16)

```



f)

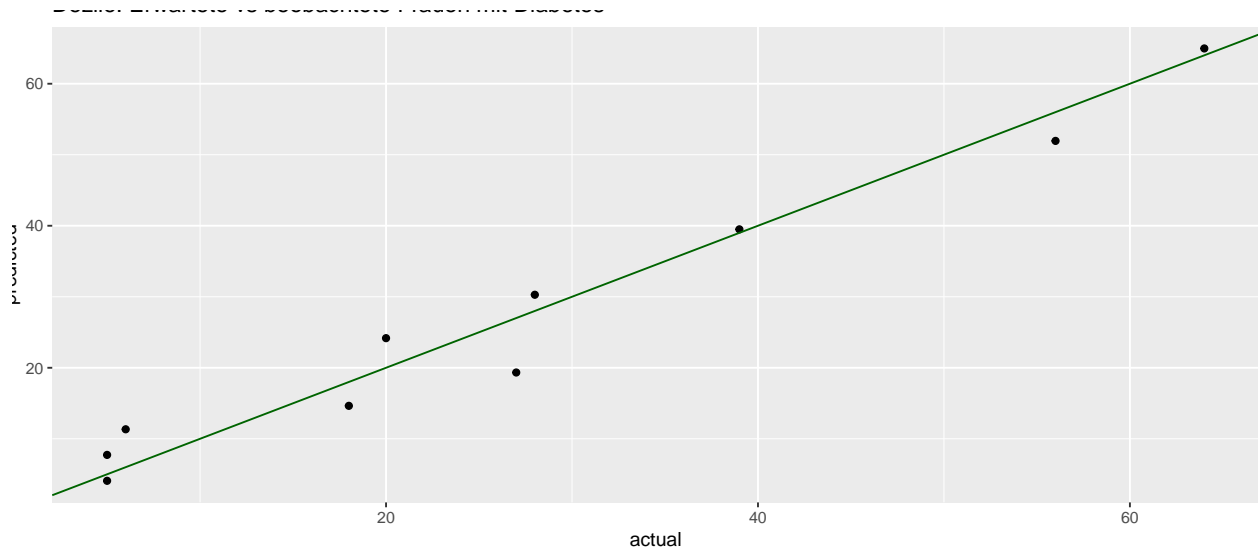
```
packageTest("StatMeasures")
packageTest("ggplot2")

diab$prob <- predict.glm(log_reg, type="response")
diab$decile <- decile(diab$prob)

actual <- c()
predicted <- c()

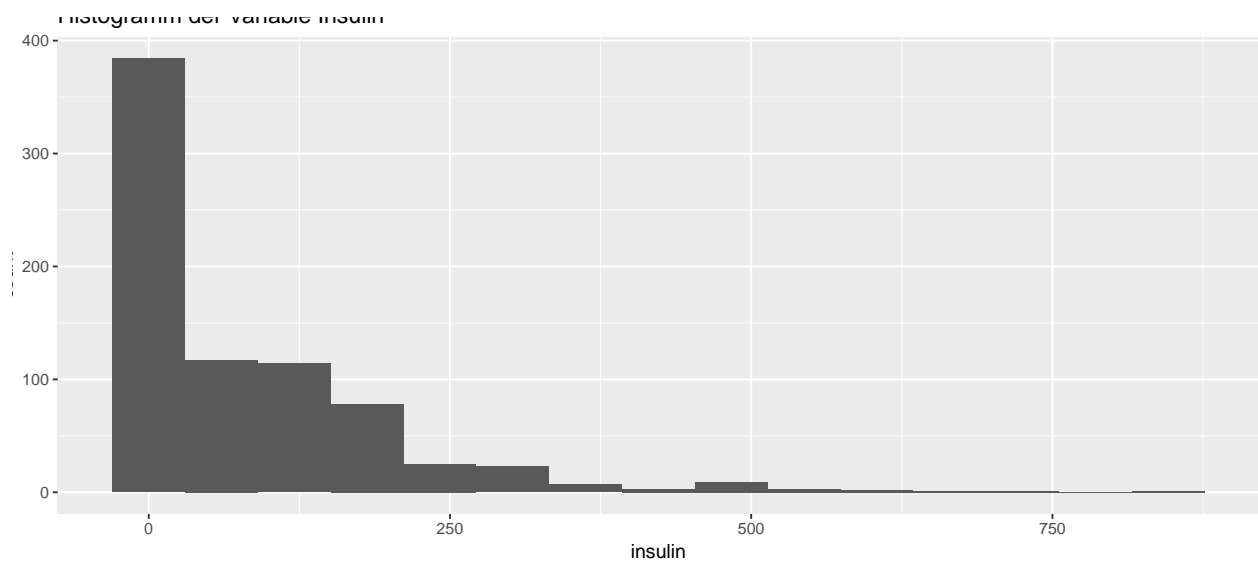
for (decile in seq(1,10)){
  predicted <- append(predicted, length(diab$pregn[diab$decile == decile])
                     *mean(diab$prob[diab$decile == decile]))
  actual <- append(actual, length(diab$pregn[(diab$decile == decile)
                                         & (diab$diabetes == 1)]))
}

decile_df <- data.frame(actual, predicted)
ggplot(aes(x=actual, y=predicted), data = decile_df) +
  geom_point() +
  geom_abline(slope=1, intercept=0, colour = 'darkgreen') +
  ggtitle("Dezile: Erwartete vs beobachtete Frauen mit Diabetes")
```

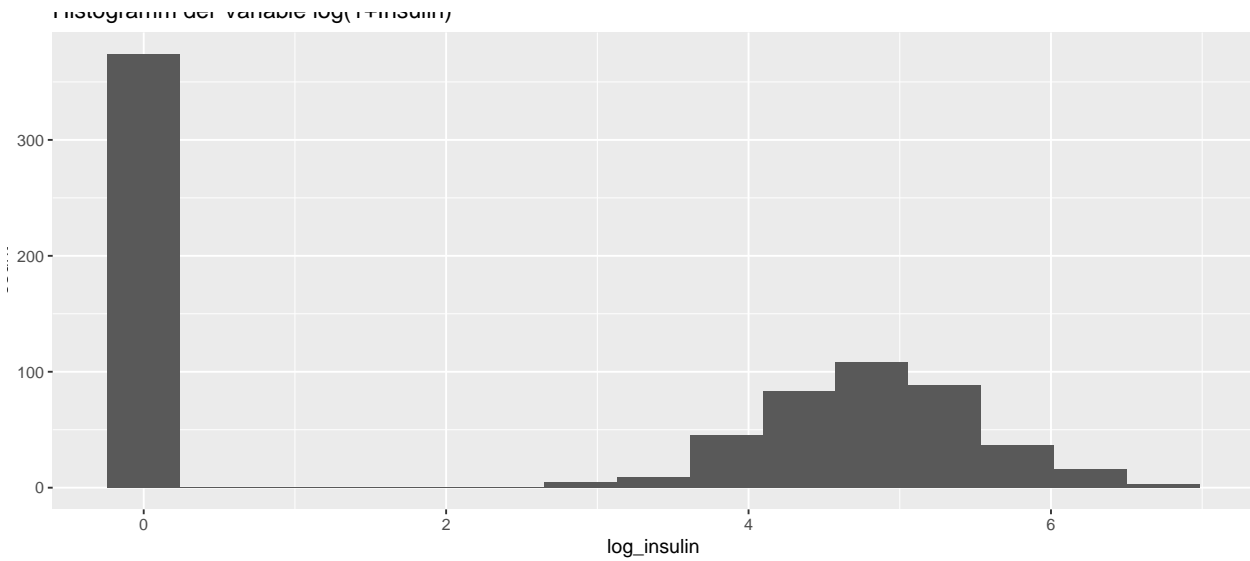


g)

```
ggplot(diab, aes(x=insulin)) +  
  geom_histogram(bins = 15) +  
  ggtitle("Histogramm der Variable Insulin")
```



```
test <- NULL  
test$log_insulin <- log1p(diab$insulin)  
test <- as.data.frame(test)  
  
ggplot(test, aes(x=log_insulin)) +  
  geom_histogram(bins = 15) +  
  ggtitle("Histogramm der Variable log(1+Insulin)")
```



```
diab$log_insulin <- log1p(diab$insulin)

log_reg1 <- glm(diabetes~pgc+bmi+log_insulin, family=binomial(link='logit'), data=diab)

test <- NULL
test$y.true <- diab$diabetes
test$y.pred <- predict.glm(log_reg1, type="response")

test$y.pred[test$y.pred > 0.5] = 1
test$y.pred[test$y.pred <= 0.5] = 0

confusion_matrix <- table(test)
confusion_matrix
#>      y.pred
#> y.true  0   1
#>    0 441  59
#>    1 129 139

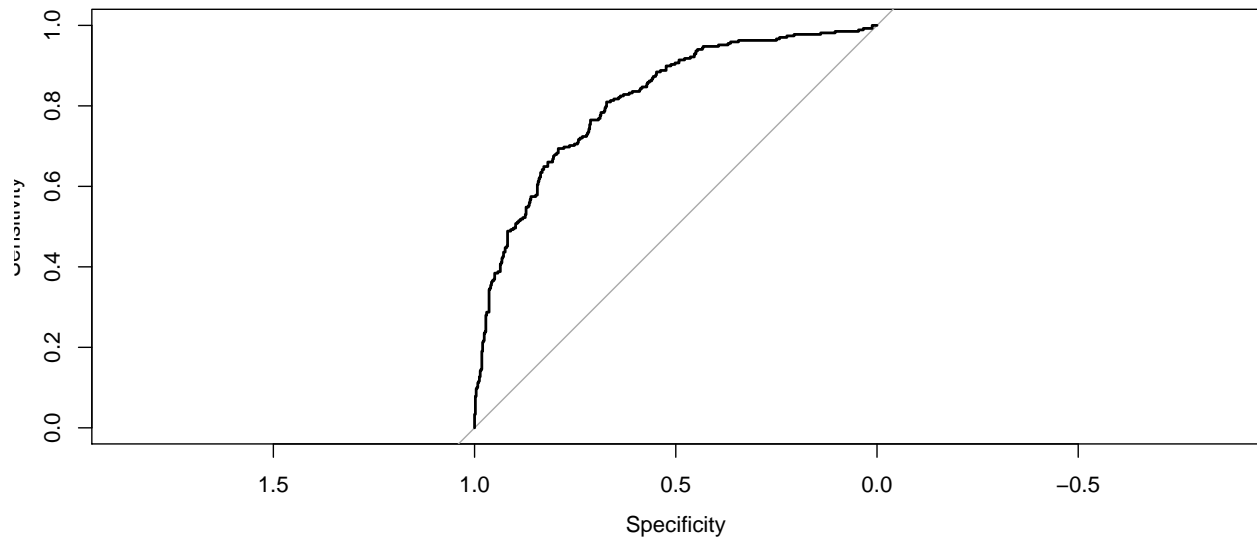
misclass_rate <- 1-sum(diag(confusion_matrix))/sum(confusion_matrix)

sensitivity <- confusion_matrix[2,2]/(confusion_matrix[2,1] + confusion_matrix[2,2])

specificity <- confusion_matrix[1,1]/(confusion_matrix[1,1] + confusion_matrix[1,2])

misclass_rate
#> [1] 0.2447917
sensitivity
#> [1] 0.5186567
specificity
#> [1] 0.882

plot.roc(diab$diabetes, predict.glm(log_reg1, type="response"))
```



Aufgabe 3

```

packageTest('class')

nor <-function(x) { (x -min(x))/(max(x)-min(x))  }
ran <- sample(1:nrow(diab), 0.5 * nrow(diab))
diab_norm <- as.data.frame(lapply(diab[,c(2,6)], nor))
#...oder alle Variablen zur Prognose?

train <- diab_norm[ran,]
test <- diab_norm[-ran,]

target_category <- diab[ran,8]
test_category <- diab[-ran,8]

knn1 <- knn(train,test,cl=target_category,k=1)
knn5 <- knn(train,test,cl=target_category,k=5)
knn10 <- knn(train,test,cl=target_category,k=10)

conf1 <- table(knn1,test_category)
conf5 <- table(knn5,test_category)
conf10 <- table(knn10,test_category)

error <- function(x){
  1 - sum(diag(x)/(sum(rowSums(x))))
}

error(conf1)
#> [1] 0.2916667
error(conf5)
#> [1] 0.2526042
error(conf10)
#> [1] 0.2604167

```