

Week 10 - Homework

STAT 420, Summer 2022, Ilya Andreev, iandre3@illinois.edu

Exercise 1 (Simulating Wald and Likelihood Ratio Tests)

In this exercise we will investigate the distributions of hypothesis tests for logistic regression. For this exercise, we will use the following predictors.

```
sample_size = 150
set.seed(120)
x1 = rnorm(n = sample_size)
x2 = rnorm(n = sample_size)
x3 = rnorm(n = sample_size)
```

Recall that

$$p(\mathbf{x}) = P[Y = 1 \mid \mathbf{X} = \mathbf{x}]$$

Consider the true model

$$\log \left(\frac{p(\mathbf{x})}{1 - p(\mathbf{x})} \right) = \beta_0 + \beta_1 x_1$$

where

- $\beta_0 = 0.4$
- $\beta_1 = -0.35$

(a) To investigate the distributions, simulate from this model 2500 times. To do so, calculate

$$P[Y = 1 \mid \mathbf{X} = \mathbf{x}]$$

for an observation, and then make a random draw from a Bernoulli distribution with that success probability. (Note that a Bernoulli distribution is a Binomial distribution with parameter $n = 1$. There is no direction function in R for a Bernoulli distribution.)

Each time, fit the model:

$$\log \left(\frac{p(\mathbf{x})}{1 - p(\mathbf{x})} \right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

Store the test statistics for two tests:

- The Wald test for $H_0 : \beta_2 = 0$, which we say follows a standard normal distribution for “large” samples
- The likelihood ratio test for $H_0 : \beta_2 = \beta_3 = 0$, which we say follows a χ^2 distribution (with some degrees of freedom) for “large” samples

```

wald = rep(0, 2500)
lrt = rep(0, 2500)

odds = exp(0.4 - 0.35 * x1)
probs = odds / (1 + odds)
for (i in 1:2500) {
  simulated_y = rbinom(length(x1), 1, probs)
  null_model = glm(simulated_y ~ x1, family=binomial)
  full_model = glm(simulated_y ~ x1 + x2 + x3, family=binomial)
  wald[i] = coef(summary(full_model))[3, 3]
  lrt[i] = anova(null_model, full_model, test = "LRT")[2, 4]
}

```

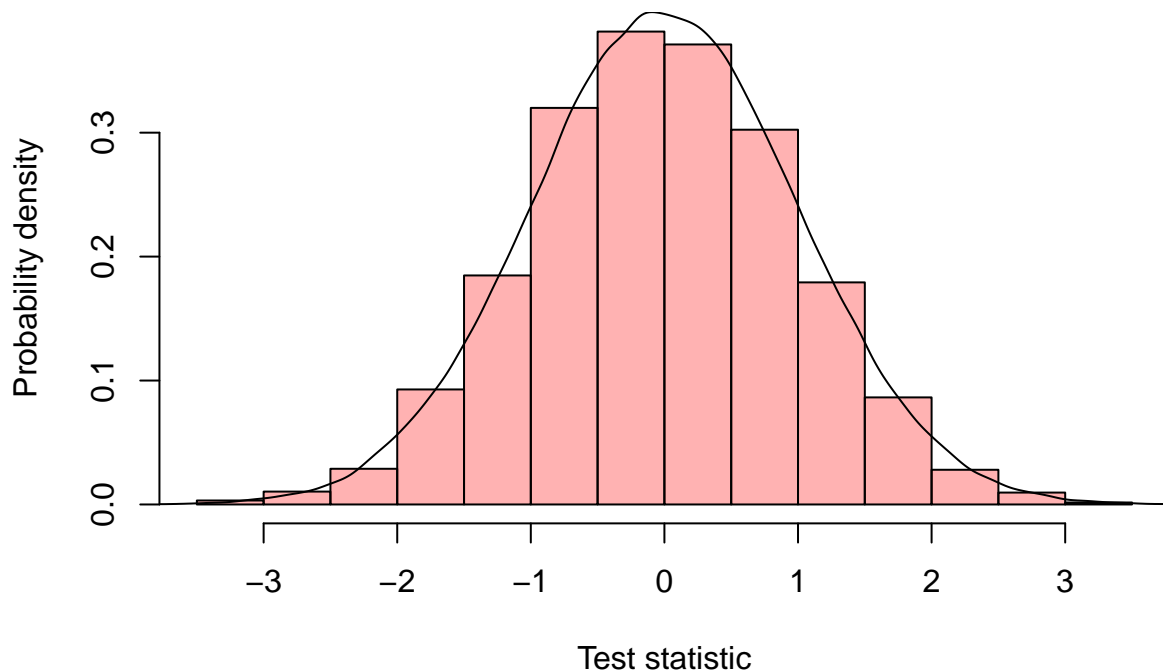
(b) Plot a histogram of the empirical values for the Wald test statistic. Overlay the density of the true distribution assuming a large sample.

```

hist(wald, freq=FALSE, main="Histogram of Wald test statistic distribution", xlab="Test statistic", ylab="Probability density", col="red", border="black")
x = rnorm(150000)
lines(density(x))

```

Histogram of Wald test statistic distribution



(c) Use the empirical results for the Wald test statistic to estimate the probability of observing a test statistic larger than 1. Also report this probability using the true distribution of the test statistic assuming a large sample.

The empirical probability is 0.1524, the true probability is 0.1587.

(d) Plot a histogram of the empirical values for the likelihood ratio test statistic. Overlay the density of the true distribution assuming a large sample.

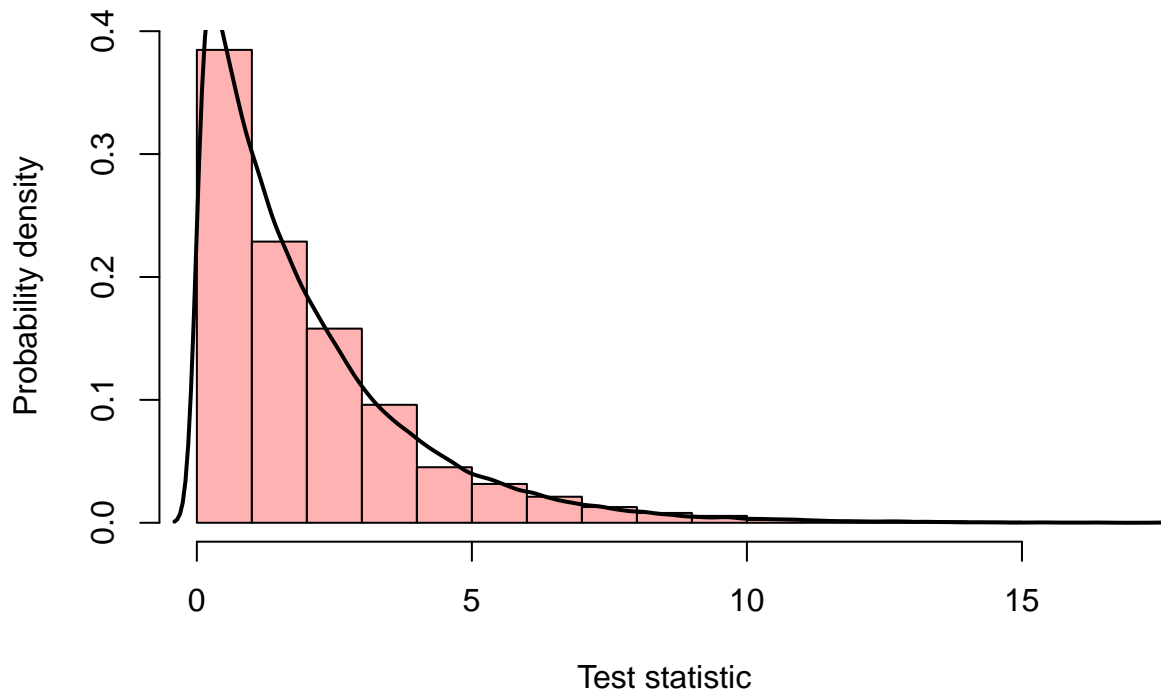
```

hist(lrt, freq=FALSE, main="Histogram of Wald test statistic distribution", xlab="Test statistic", ylab="Probability density", col="red", border="black")
x = rchisq(150000, 2)

```

```
lines(density(x), lw=2)
```

Histogram of Wald test statistic distribution



(e) Use the empirical results for the likelihood ratio test statistic to estimate the probability of observing a test statistic larger than 5. Also report this probability using the true distribution of the test statistic assuming a large sample.

The empirical probability is 0.0872, the true probability is 0.0821.

(f) Repeat (a)-(e) but with simulation using a smaller sample size of 10. Based on these results, is this sample size large enough to use the standard normal and χ^2 distributions in this situation? Explain.

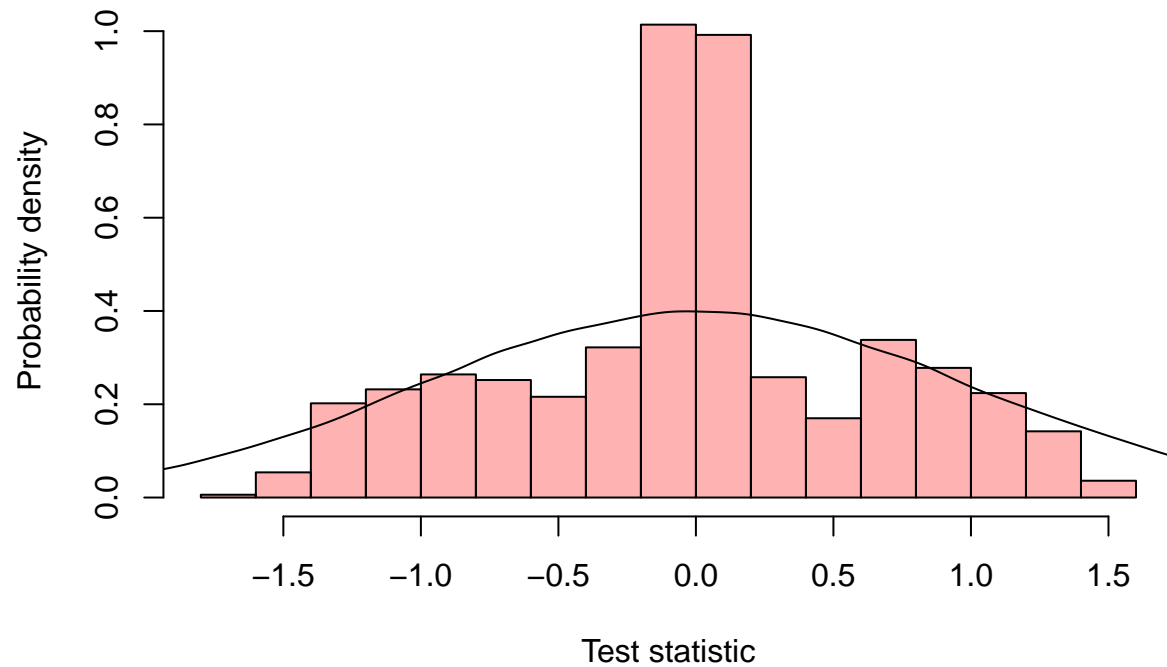
```
options(warn=-1)
sample_size = 10
set.seed(120)
x1 = rnorm(n = sample_size)
x2 = rnorm(n = sample_size)
x3 = rnorm(n = sample_size)

wald = rep(0, 2500)
lrt = rep(0, 2500)

odds = exp(0.4 - 0.35 * x1)
probs = odds / (1 + odds)
for (i in 1:2500) {
  simulated_y = rbinom(length(x1), 1, probs)
  null_model = glm(simulated_y ~ x1, family=binomial)
  full_model = glm(simulated_y ~ x1 + x2 + x3, family=binomial)
  wald[i] = coef(summary(full_model))[3, 3]
  lrt[i] = anova(null_model, full_model, test = "LRT")[2, 4]
}
```

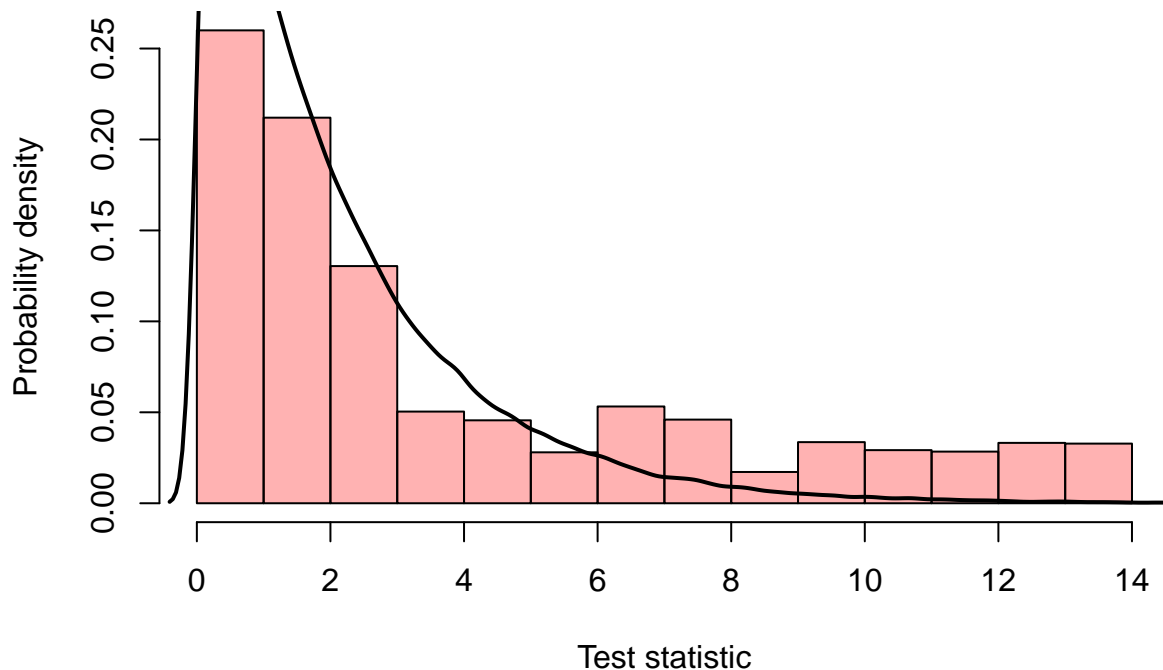
```
hist(wald, freq=FALSE, main="Histogram of Wald test statistic distribution", xlab="Test statistic", ylab="Probability density")
x = rnorm(150000)
lines(density(x))
```

Histogram of Wald test statistic distribution



```
hist(lrt, freq=FALSE, main="Histogram of Wald test statistic distribution", xlab="Test statistic", ylab="Probability density")
x = rchisq(150000, 2)
lines(density(x), lw=2)
```

Histogram of Wald test statistic distribution



For the Wald test, the empirical probability is 0.0804, the true probability is 0.1587. For the likelihood ratio test, the empirical probability is 0.3016, the true probability is 0.0821.

As can be seen, there is a substantial difference in the empirical Wald and LRT test statistic distributions when the sample size is small. The conclusion is that the sample size of 10 is not large enough.

Exercise 2 (Surviving the Titanic)

For this exercise use the `ptitanic` data from the `rpart.plot` package. (The `rpart.plot` package depends on the `rpart` package.) Use `?rpart.plot::ptitanic` to learn about this dataset. We will use logistic regression to help predict which passengers aboard the [Titanic](#) will survive based on various attributes.

```
# install.packages("rpart")
# install.packages("rpart.plot")
library(rpart)
library(rpart.plot)
data("ptitanic")
```

For simplicity, we will remove any observations with missing data. Additionally, we will create a test and train dataset.

```
ptitanic = na.omit(ptitanic)
set.seed(2021)
trn_idx = sample(nrow(ptitanic), 300)
ptitanic_trn = ptitanic[trn_idx, ]
ptitanic_tst = ptitanic[-trn_idx, ]
```

(a) Consider the model

$$\log \left(\frac{p(\mathbf{x})}{1 - p(\mathbf{x})} \right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_3 x_4$$

where

$$p(\mathbf{x}) = P[Y = 1 \mid \mathbf{X} = \mathbf{x}]$$

is the probability that a certain passenger survives given their attributes and

- x_1 is a dummy variable that takes the value 1 if a passenger was 2nd class.
- x_2 is a dummy variable that takes the value 1 if a passenger was 3rd class.
- x_3 is a dummy variable that takes the value 1 if a passenger was male.
- x_4 is the age in years of a passenger.

Fit this model to the training data and report its deviance.

```
ptitanic_trn$pclass = as.factor(ptitanic_trn$pclass)
ptitanic_trn$sex = as.factor(ptitanic_trn$sex)
ptitanic_trn$survived = as.factor(ptitanic_trn$survived)
ptitanic_tst$pclass = as.factor(ptitanic_tst$pclass)
ptitanic_tst$sex = as.factor(ptitanic_tst$sex)
ptitanic_tst$survived = as.factor(ptitanic_tst$survived)

model = glm(survived ~ pclass + sex + age + age * sex, ptitanic_trn, family=binomial)
summary(model)
```

```
##
## Call:
## glm(formula = survived ~ pclass + sex + age + age * sex, family = binomial,
##      data = ptitanic_trn)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.400   -0.665   -0.364    0.688    2.588
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  2.795764   0.730785   3.83 0.00013 ***
## pclass2nd    -1.496417   0.490649  -3.05 0.00229 **
## pclass3rd    -2.836925   0.486207  -5.83 5.4e-09 ***
## sexmale      -0.473149   0.697951  -0.68 0.49783
## age           0.000547   0.018387   0.03 0.97625
## sexmale:age  -0.077239   0.024128  -3.20 0.00137 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 396.42  on 299  degrees of freedom
## Residual deviance: 259.17  on 294  degrees of freedom
## AIC: 271.2
##
## Number of Fisher Scoring iterations: 5
```

The deviance of the model is 396.4236.

(b) Use the model fit in (a) and an appropriate statistical test to determine if class played a significant role in surviving on the Titanic. Use $\alpha = 0.01$. Report:

- The null hypothesis of the test
- The test statistic of the test
- The p-value of the test
- A statistical decision
- A practical conclusion

The null hypothesis:

$$H_0 : \beta_1 = \beta_2 = 0$$

```
small_model = glm(survived ~ sex + age + age * sex, ptitanic_trn, family=binomial)
res = anova(small_model, model, test="LRT")
res
```

```
## Analysis of Deviance Table
##
## Model 1: survived ~ sex + age + age * sex
## Model 2: survived ~ pclass + sex + age + age * sex
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1         296         304
## 2         294         259  2      45.1  1.6e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Test statistic: 45.0617. P-value: 1.6405×10^{-10} . Statistical decision: reject the null hypothesis. Practical conclusion: class played a significant role in surviving on Titanic.

(c) Use the model fit in (a) and an appropriate statistical test to determine if an interaction between age and sex played a significant role in surviving on the Titanic. Use $\alpha = 0.01$. Report:

- The null hypothesis of the test
- The test statistic of the test
- The p-value of the test
- A statistical decision
- A practical conclusion

The null hypothesis:

$$H_0 : \beta_5 = 0$$

```
small_model = glm(survived ~ pclass + sex + age, ptitanic_trn, family=binomial)
res = anova(small_model, model, test="LRT")
res
```

```
## Analysis of Deviance Table
##
## Model 1: survived ~ pclass + sex + age
## Model 2: survived ~ pclass + sex + age + age * sex
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1         295         270
## 2         294         259  1      11.4  0.00075 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Test statistic: 11.3717. P-value: 0.0007. Statistical decision: reject the null hypothesis. Practical conclusion: interaction between age and sex played a significant role in surviving on Titanic.

(d) Use the model fit in (a) as a classifier that seeks to minimize the misclassification rate. Classify each of the passengers in the test dataset. Report the misclassification rate, the sensitivity, and the specificity of this classifier. (Use survived as the positive class.)

```
predictions = ifelse(predict(model, ptitanic_tst, type="response") > 0.5, "survived", "died") != ptitanic_tst$survived

make_conf_mat = function(predicted, actual) {
  table(predicted = predicted, actual = actual)
}

conf_mat = make_conf_mat(predictions, ptitanic_tst$survived)
conf_mat

##           actual
## predicted died survived
##    FALSE  404      180
##     TRUE   27      135

get_sens = function(conf_mat) {
  conf_mat[2, 2] / sum(conf_mat[, 2])
}

get_spec = function(conf_mat) {
  conf_mat[1, 1] / sum(conf_mat[, 1])
}
```

The misclassification rate is 0.2172. The sensitivity is 0.4286. The specificity is 0.9374.

Exercise 3 (Breast Cancer Detection)

For this exercise we will use data found in `wisc-train.csv` and `wisc-test.csv`, which contain train and test data, respectively. `wisc.csv` is provided but not used. This is a modification of the Breast Cancer Wisconsin (Diagnostic) dataset from the UCI Machine Learning Repository. Only the first 10 feature variables have been provided. (And these are all you should use.)

- [UCI Page](#)
- [Data Detail](#)

You should consider coercing the response to be a factor variable if it is not stored as one after importing the data.

(a) The response variable `class` has two levels: M if a tumor is malignant, and B if a tumor is benign. Fit three models to the training data.

- An additive model that uses `radius`, `smoothness`, and `texture` as predictors
- An additive model that uses all available predictors
- A model chosen via backwards selection using AIC. Use a model that considers all available predictors as well as their two-way interactions for the start of the search.

For each, obtain a 5-fold cross-validated misclassification rate using the model as a classifier that seeks to minimize the misclassification rate. Based on this, which model is best? Relative to the best, are the other two underfitting or over fitting? Report the test misclassification rate for the model you picked as the best.

```
options(warn=-1)
wisc_train = read.csv("wisc-train.csv")
wisc_train$class = as.factor(wisc_train$class)
wisc_test = read.csv("wisc-test.csv")
```



```

wisc_test$class = as.factor(wisc_test$class)

m1 = glm(class ~ radius + smoothness + texture, wisc_train, family=binomial)
m2 = glm(class ~ ., wisc_train, family=binomial)
m3 = step(glm(class ~ . ^ 2, wisc_train, family=binomial), trace=0)

options(warn=-1)
library(boot)
cv.glm(wisc_train, m1, K = 5)$delta[1]

## [1] 0.08047

cv.glm(wisc_train, m2, K = 5)$delta[1]

## [1] 0.1041

cv.glm(wisc_train, m3, K = 5)$delta[1]

## [1] 0.0909

predictions1 = ifelse(predict(m1, wisc_test, type="response") > 0.5, "M", "B")
mean(predictions1 != wisc_test$class)

## [1] 0.08955

predictions2 = ifelse(predict(m2, wisc_test, type="response") > 0.5, "M", "B")
mean(predictions2 != wisc_test$class)

## [1] 0.1173

predictions3 = ifelse(predict(m3, wisc_test, type="response") > 0.5, "M", "B")
mean(predictions3 != wisc_test$class)

## [1] 0.1514

```

The first, the smallest model, has the lowest 5-fold cross-validation misclassification rate. As seen by the misclassification rate on the test dataset, the other two larger models are overfitting.

(b) In this situation, simply minimizing misclassifications might be a bad goal since false positives and false negatives carry very different consequences. Consider the M class as the “positive” label. Consider each of the probabilities stored in `cutoffs` in the creation of a classifier using the **additive** model fit in (a).

```
cutoffs = seq(0.01, 0.99, by = 0.01)
```

That is, consider each of the values stored in `cutoffs` as c . Obtain the sensitivity and specificity in the test set for each of these classifiers. Using a single graphic, plot both sensitivity and specificity as a function of the cutoff used to create the classifier. Based on this plot, which cutoff would you use? (0 and 1 have not been considered for coding simplicity. If you like, you can instead consider these two values.)

$$\hat{C}(\mathbf{x}) = \begin{cases} 1 & \hat{p}(\mathbf{x}) > c \\ 0 & \hat{p}(\mathbf{x}) \leq c \end{cases}$$

```

sensitivities = rep(0, length(cutoffs))
specificities = rep(0, length(cutoffs))
i = 1
for (c in cutoffs) {
  predictions = ifelse(predict(m1, wisc_test, type="response") > c, "M", "B")
  conf_mat = make_conf_mat(predictions, wisc_test$class)
  sensitivities[i] = get_sens(conf_mat)
}

```

```

specificities[i] = get_spec(conf_mat)
i = i + 1
}

plot(cutoffs, sensitivities, col="blue", type="l", ylim=c(0, 1), xlim=c(0, 1), main="Sensitivity vs Specificity",
     par(new=TRUE))
plot(cutoffs, specificities, col="red", type="l", ylim=c(0, 1), xlim=c(0, 1), xlab="Cutoff", ylab="Metric")
legend(0.1, 0.5, legend=c("Sensitivity", "Specificity"),
      col=c("blue", "red"), lty=1, cex=0.8)

```

Sensitivity vs Specificity at different decision boundaries

