

# Predicting cancer diagnosis from cell imaging data

Nathan Anderson and Adam Koller

4/7/2022

## Abstract

This report investigates tumor diagnosis (malignant or benign) as a function of various cell nuclei measurements from a fine needle aspirate image of a breast tumor mass. The measurements include average cell nuclei radius, texture, perimeter, area, smoothness, compactness, concavity, concave points, symmetry, and fractal dimension. The primary goal of the report is to fit a multiple logistic regression model using the aforementioned predictors that can predict tumor diagnosis with reasonable accuracy and precision. Additionally, we would like to determine which predictors are most predictive of tumor diagnosis. From our best model, which only included a subset of the predictors, we found that mean texture, smoothness, square root concavity, log compactness, and log area are all significant predictors of tumor diagnosis. We found five significant interactions between these five predictor variables. Altogether, 96.2% of malignant cases and 96.9% of benign cases were correctly predicted by our final model.

## Introduction

Each year, approximately 250,000 people are diagnosed with breast cancer in the United States alone. To treat these patients, the US spends upwards of 20 billion dollars each year. Early detection of cancer through the assessment of tumors as benign or malignant is essential to limiting morbidities and health care costs associated with breast cancer. Obtaining an accurate assessment of tumor diagnosis is of great importance to doctors, patients, and their families. One method of tumor diagnosis is by analyzing cell features and measurements obtained by a fine needle aspirate (FNA) biopsy procedure. The images generated by the FNA are analyzed by a cytotechnologist or pathologist for signs of malignancy. Importantly, multiple logistic regression models could be valuable resources that expedite the process and help inform the doctors making the final diagnosis.

This report examines data from 569 breast tumor samples from patients in Wisconsin. In addition to our response variable (diagnosis), there are various cell nuclei measurements commonly calculated by computers from FNA images. The predictor variables included in our report are:

- a) mean radius (mean of distances from center to points on the perimeter)
- b) mean texture (standard deviation of gray-scale values)
- c) mean perimeter
- d) mean area
- e) mean smoothness (local variation in radius lengths)
- f) mean compactness ( $\text{perimeter}^2 / \text{area} - 1$ )
- g) mean concavity (severity of concave portions of the contour)
- h) mean concave points (number of concave portions of the contour)
- i) mean symmetry
- j) mean fractal dimension (“coastline approximation” - 1)

All predictor variables are quantitative.

The primary goal of this report is to determine what cell nuclei measurements are most predictive of tumor

diagnosis. Additionally, we want to assess how effective a combination of cell nuclei measurements can be at predicting tumor diagnosis.

## Loading the data

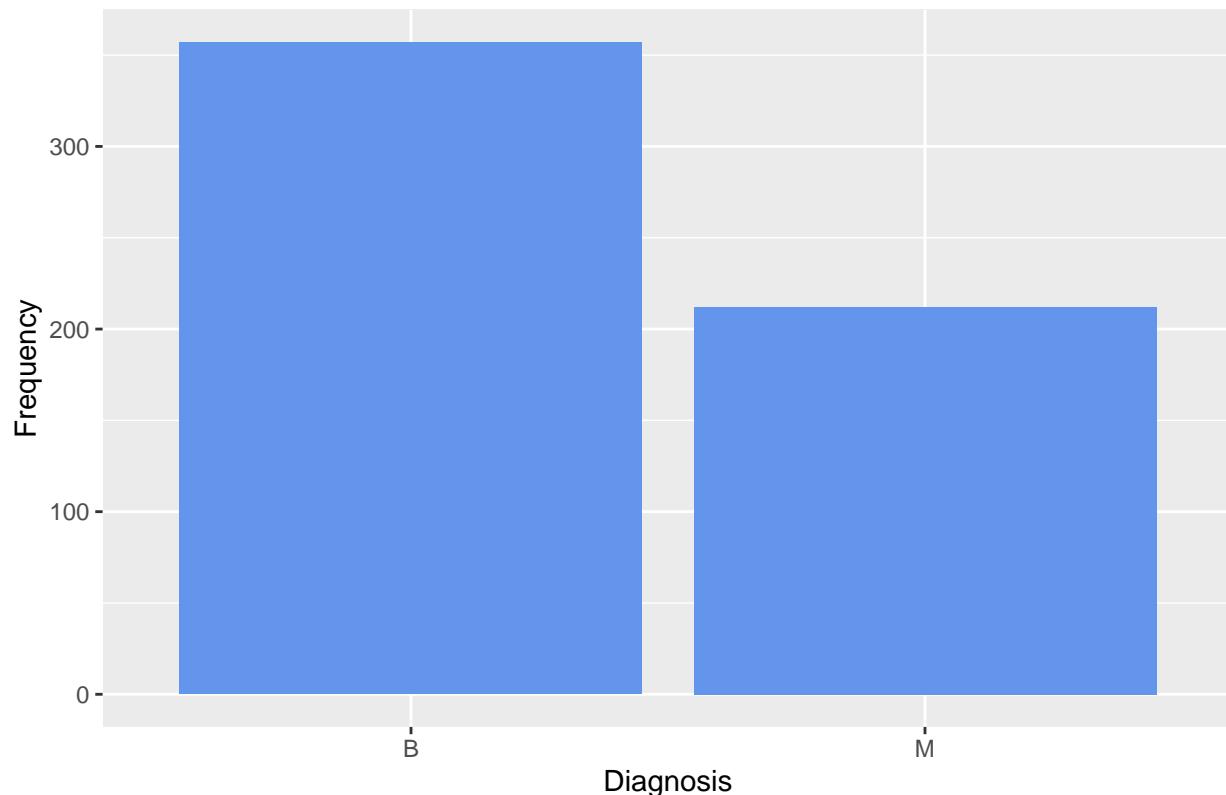
```
library(readxl)
wisconsin <- read_excel("G:/Shared drives/math 327 project/wisconsin_cancer.xlsx")
#wisconsin <- read_excel("/Users/nathananderson/Downloads/wisconsin_cancer.xlsx")
wisconsin$diagnosis <- as.factor(wisconsin$diagnosis)
wisconsin <- wisconsin[2:12]
diagnosis = wisconsin$diagnosis
radius_mean = wisconsin$radius_mean
texture_mean = wisconsin$texture_mean
perimeter_mean = wisconsin$perimeter_mean
area_mean = wisconsin$area_mean
smoothness_mean = wisconsin$smoothness_mean
compactness_mean = wisconsin$compactness_mean
concavity_mean = wisconsin$concavity_mean
concave_points_mean = wisconsin$`concave points_mean`
symmetry_mean = wisconsin$symmetry_mean
fractal_dimension_mean = wisconsin$fractal_dimension_mean
```

## Data Characteristics

```
#Response variable characteristics
library(ggplot2)

## Warning: package 'ggplot2' was built under R version 4.0.5
ggplot(wisconsin, aes(x = diagnosis)) + geom_bar(fill="cornflowerblue") + labs(x = "Diagnosis", y = "Frequency", title = "Frequency of Benign (B) and Malignant (M) Tumors")
```

## Frequency of Benign (B) and Malignant (M) diagnoses



```
summary(wisconsin$diagnosis)
```

```
##   B   M
## 357 212
```

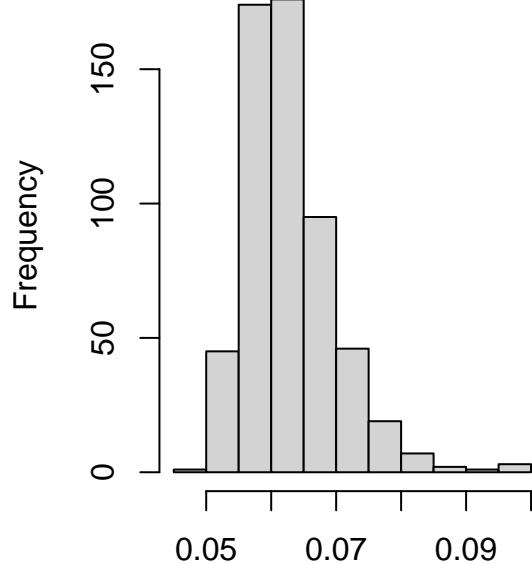
Our response variable, diagnosis, shows that 357 patient tumors were classified as benign and 212 patient tumors were classified as malignant. Our overall sample size is 569.

## Predictor variable distributions

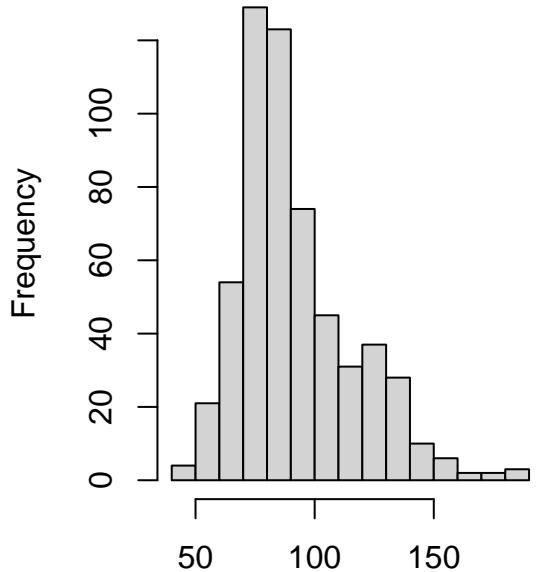
```
par(mfrow = c(1, 2))
hist(wisconsin$fractal_dimension_mean)

hist(wisconsin$perimeter_mean)
```

ogram of wisconsin\$fractal\_dimension histogram of wisconsin\$perimeter\_



wisconsin\$fractal\_dimension\_mean

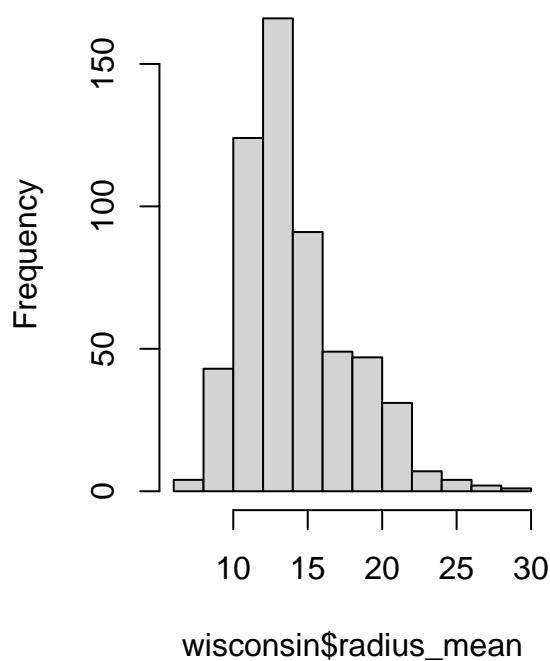


wisconsin\$perimeter\_mean

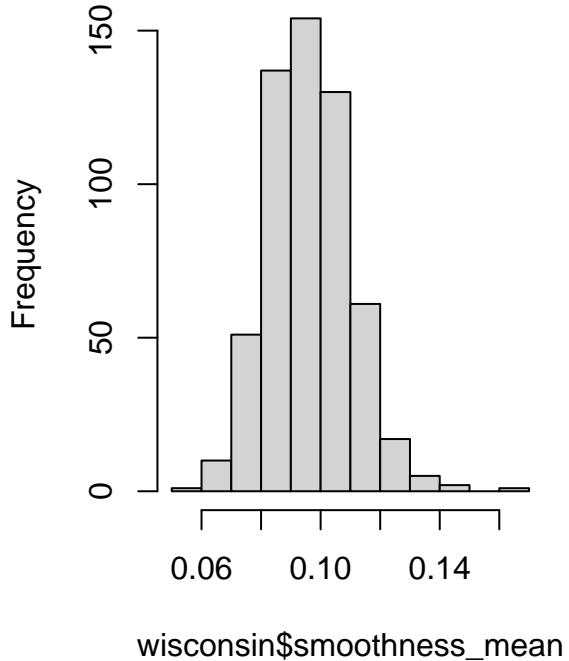
```
hist(wisconsin$radius_mean)
```

```
hist(wisconsin$smoothness_mean)
```

## Histogram of wisconsin\$radius\_mean histogram of wisconsin\$smoothness\_mean



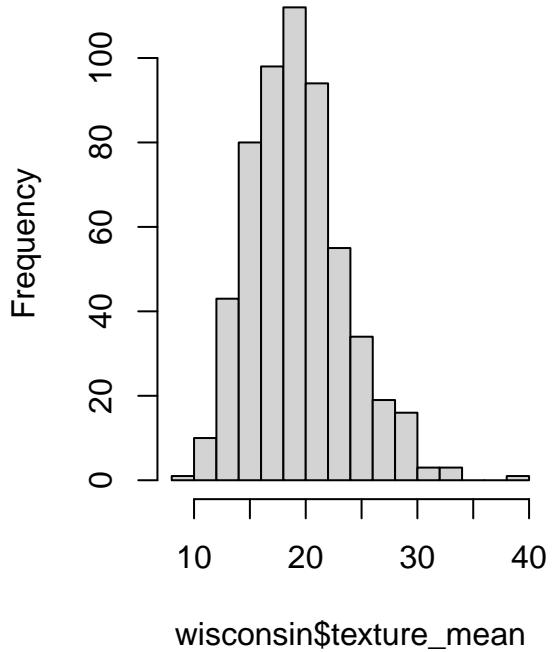
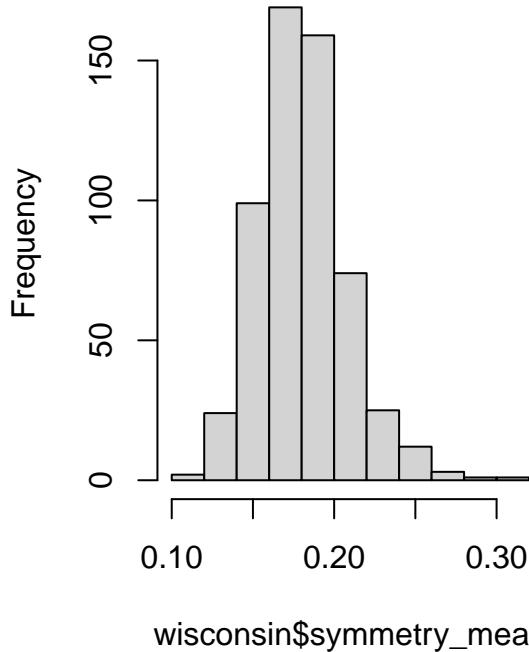
wisconsin\$radius\_mean



wisconsin\$smoothness\_mean

```
hist(wisconsin$symmetry_mean)  
hist(wisconsin$texture_mean)
```

## listogram of wisconsin\$symmetry\_ Histogram of wisconsin\$texture\_m

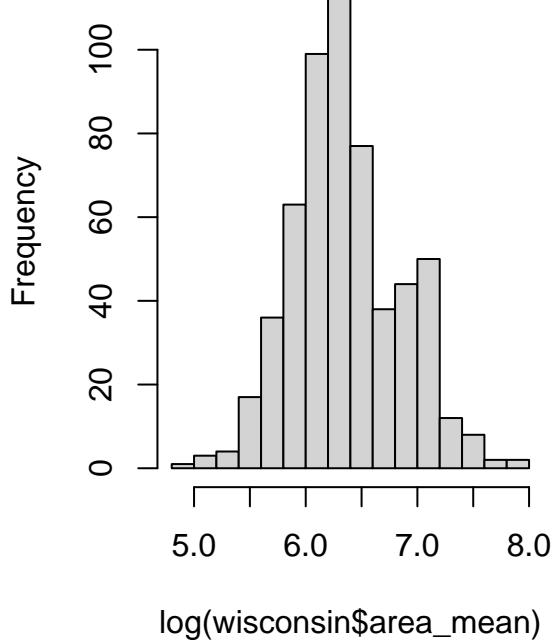
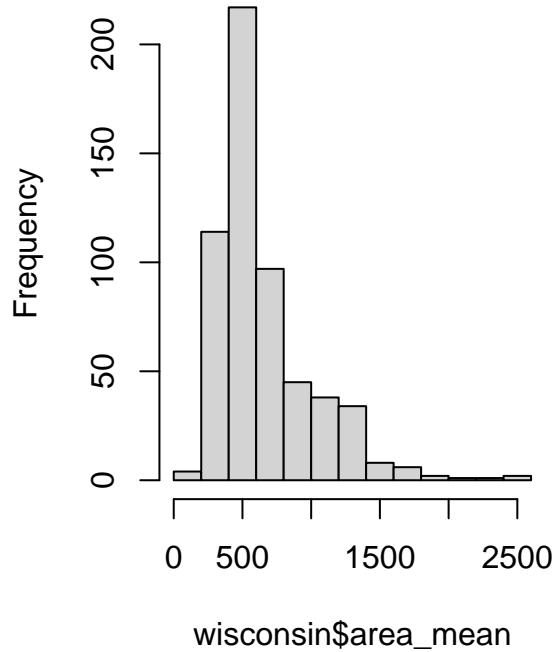


Mean fractal dimension, perimeter, radius, smoothness, symmetry, and texture are all approximately normally distributed or only mildly right skewed. As a result, we will not be transforming these predictor variables.

```
par(mfrow = c(1,2))
hist(wisconsin$area_mean)
hist(log(wisconsin$area_mean))
```

## Histogram of wisconsin\$area\_mean

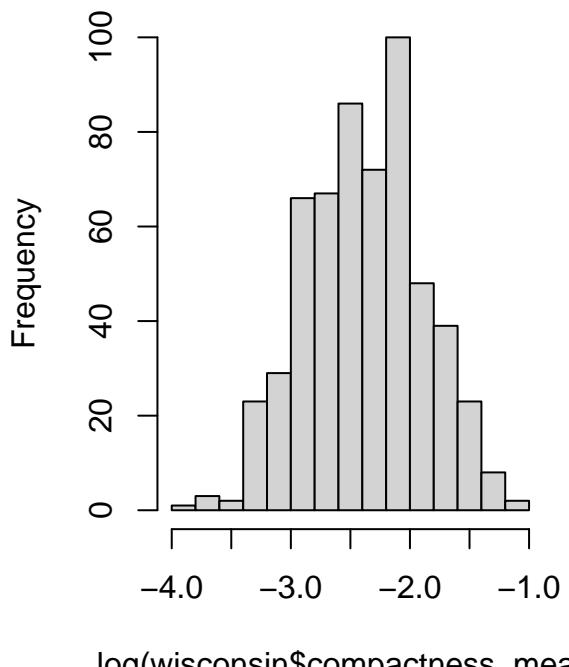
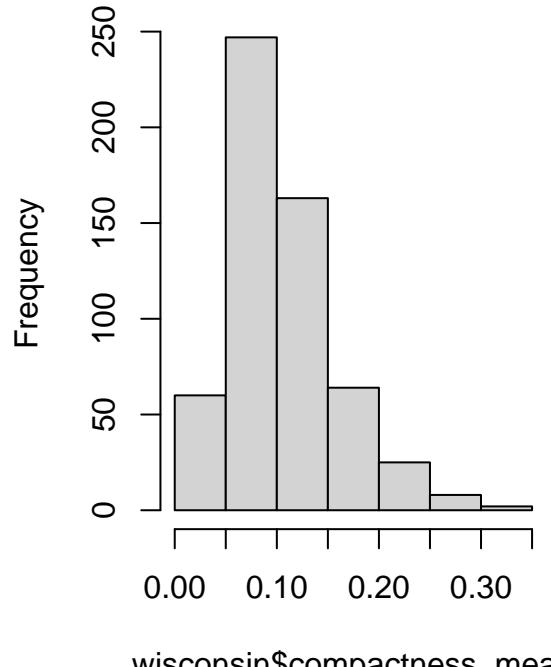
## Histogram of log(wisconsin\$area\_mean)



Mean area was not strongly right skewed. A log transformation of mean area resulted in a more symmetrical and normal distribution and will be proceeded with for further analysis.

```
par(mfrow = c(1,2))
hist(wisconsin$compactness_mean)
hist(log(wisconsin$compactness_mean))
```

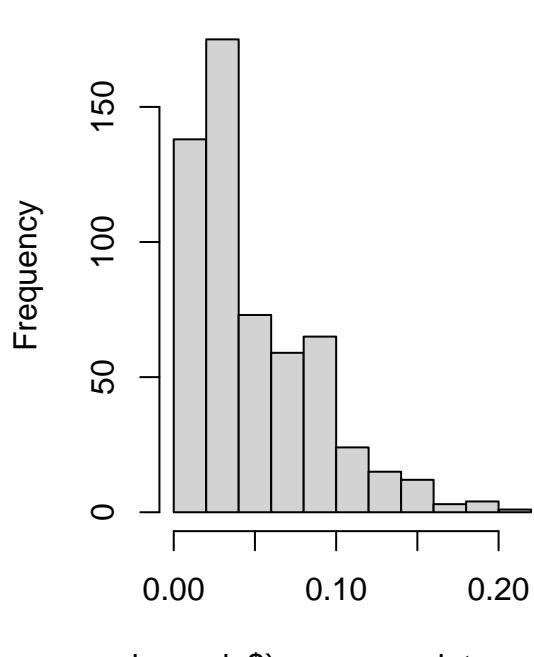
## histogram of wisconsin\$compactness



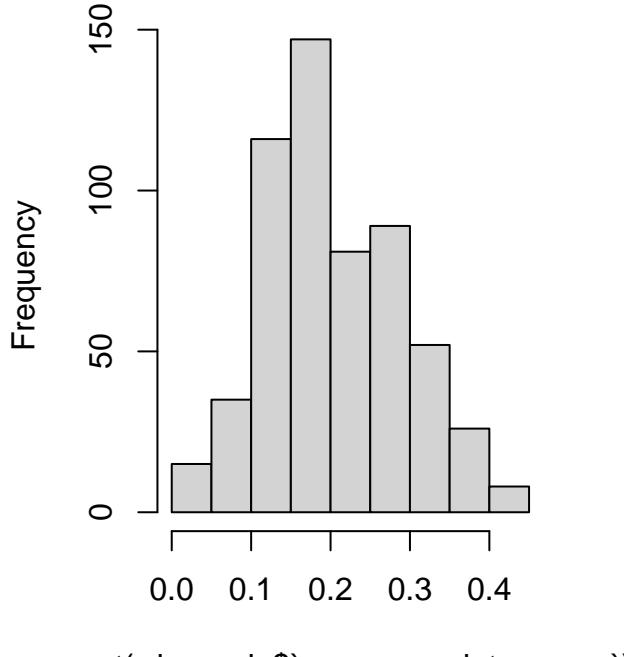
Mean compactness was strongly right skewed. A log transformation of mean compactness resulted in a more symmetrical and normal distribution and will be proceeded with in further analysis.

```
par(mfrow = c(1, 2))
hist(wisconsin$`concave points_mean`)
hist(sqrt(wisconsin$`concave points_mean`))
```

## ogram of wisconsin\$`concave pointam of sqrt(wisconsin\$`concave poi



wisconsin\$`concave points\_mean`

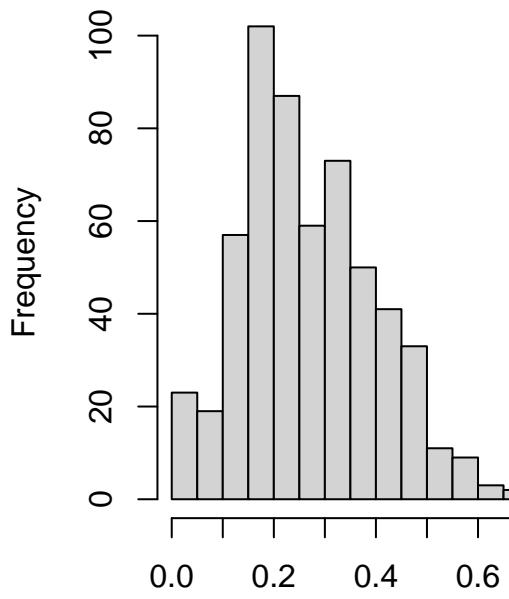
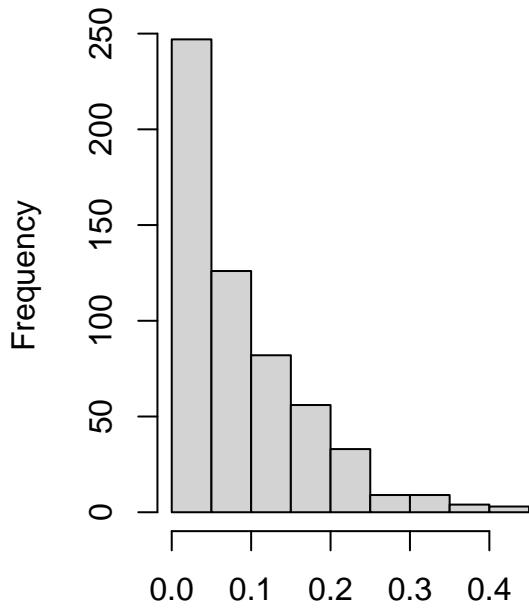


sqrt(wisconsin\$`concave points\_mean`)

Mean concave points was severely skewed. A square root transformations resulted in a more symmetrical and normal distribution and will be proceeded with in further analysis.

```
par(mfrow = c(1, 2))
hist(wisconsin$concavity_mean)
hist(sqrt(wisconsin$concavity_mean))
```

**histogram of wisconsin\$concavity\_mean**



Mean concavity is strongly skewed. A square root transformation resulted in a more symmetrical and normal distribution and will be proceeded with in further analysis.

```
#Adding transformations of area, compactness, concavity, and concave points
wisconsin$logarea_mean <- log(area_mean)
wisconsin$logcompactness_mean <- log(compactness_mean)
wisconsin$sqrtconcavity_mean <- sqrt(concavity_mean)
wisconsin$sqrtconcave_points_mean <- sqrt(concave_points_mean)
```

Histograms for the final set of predictor variables:

```
# Histograms of predictor variables
library(ggplot2)
library(tidyr)

## Warning: package 'tidyr' was built under R version 4.0.4
library(dplyr)

## Warning: package 'dplyr' was built under R version 4.0.4
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##     filter, lag

## The following objects are masked from 'package:base':
##
```

```

##      intersect, setdiff, setequal, union
plong = wisconsin[,c(1:4, 6, 10:15)] %>% group_by(diagnosis) %>%
  pivot_longer (2:11)

ggplot(plong, aes(value)) +
  geom_histogram(bins = 16, aes(fill=diagnosis, color=diagnosis), alpha=0.5, position = "identity") +
  facet_wrap(~name, scales = 'free_x')

```

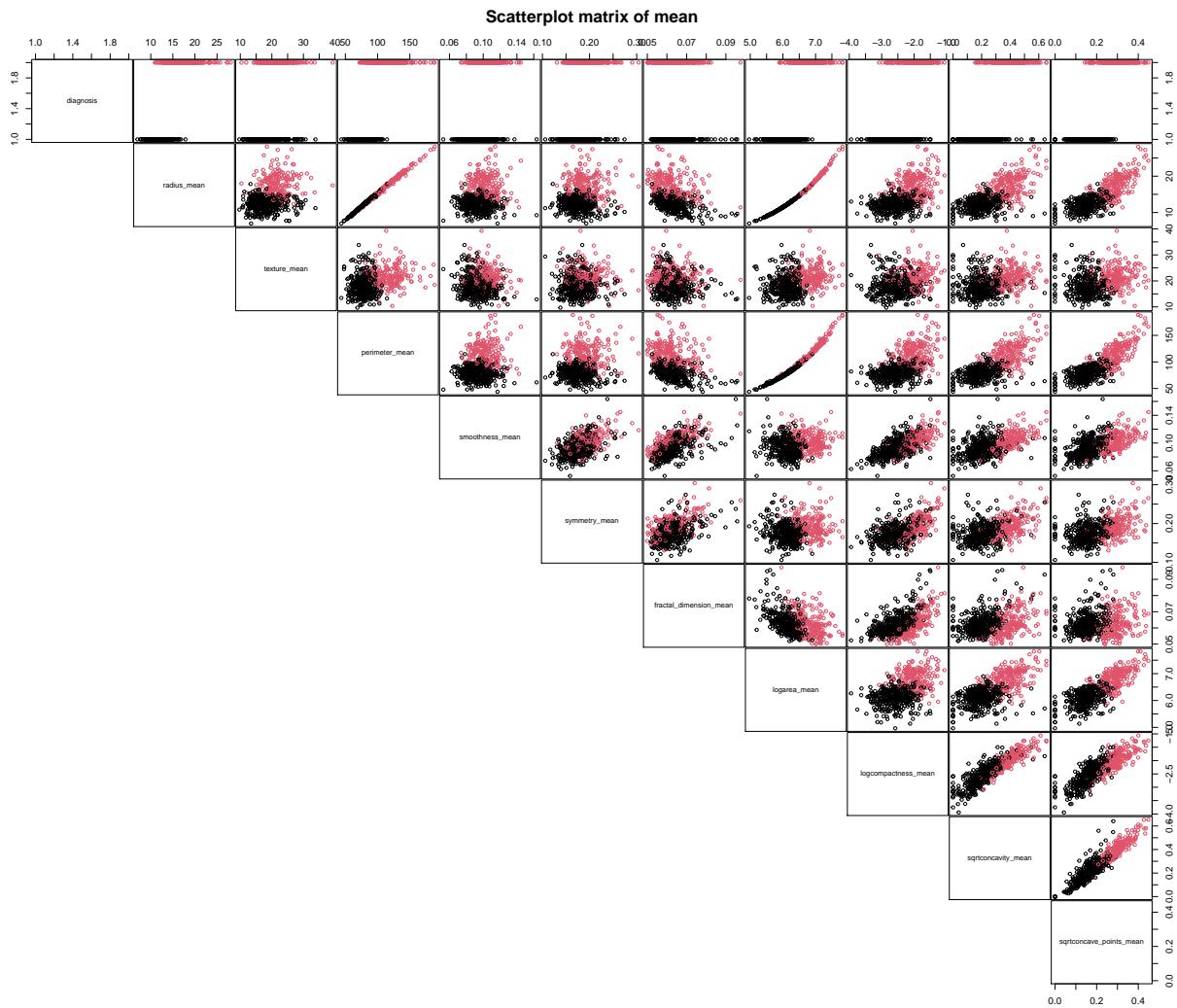


A two way frequency table categorized by diagnosis shows that multiple predictor variables, particularly mean perimeter, log area, radius, log compactness, sqrt concave points, and sqrt concavity, are associated with diagnosis. Mean fractal dimension and symmetry appear as though they will be less predictive of diagnosis.

```

pairs(wisconsin[, c(1:4, 6, 10:15)], col=diagnosis, pch=1, cex = 0.75, lower.panel = NULL, gap = 0.1, m)

```



```
cormat= cor(wisconsin[, c(2:4, 6, 10:15)])
round(cormat,3)
```

	radius_mean	texture_mean	perimeter_mean	smoothness_mean
## radius_mean	1.000	0.324	0.998	0.171
## texture_mean	0.324	1.000	0.330	-0.023
## perimeter_mean	0.998	0.330	1.000	0.207
## smoothness_mean	0.171	-0.023	0.207	1.000
## symmetry_mean	0.148	0.071	0.183	0.558
## fractal_dimension_mean	-0.312	-0.076	-0.261	0.585
## logarea_mean	0.986	0.320	0.982	0.141
## logcompactness_mean	0.502	0.222	0.548	0.683
## sqrtconcavity_mean	0.680	0.294	0.716	0.527
## sqrtconcave_points_mean	0.798	0.268	0.824	0.563
## symmetry_mean	0.148	-0.312	0.986	
## radius_mean	0.071	-0.076	0.320	
## perimeter_mean	0.183	-0.261	0.982	
## smoothness_mean	0.558	0.585	0.141	
## symmetry_mean	1.000	0.480	0.120	

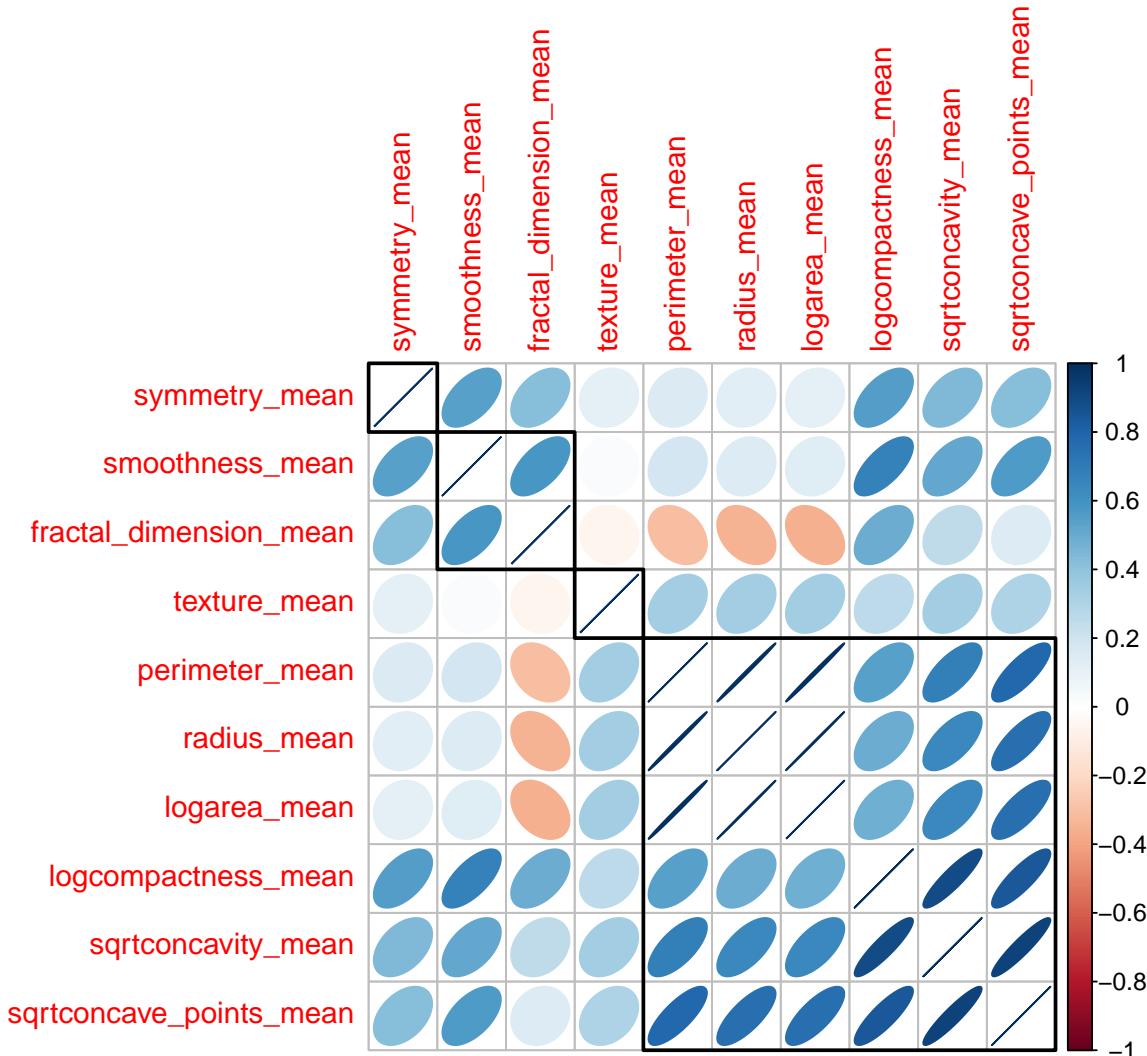
```

## fractal_dimension_mean      0.480      1.000     -0.356
## logarea_mean                0.120     -0.356      1.000
## logcompactness_mean         0.562      0.525     0.481
## sqrtconcavity_mean          0.468      0.307     0.661
## sqrtconcave_points_mean    0.438      0.158     0.787
##                               logcompactness_mean sqrtconcavity_mean
## radius_mean                  0.502      0.680
## texture_mean                 0.222      0.294
## perimeter_mean                0.548      0.716
## smoothness_mean               0.683      0.527
## symmetry_mean                 0.562      0.468
## fractal_dimension_mean        0.525      0.307
## logarea_mean                  0.481      0.661
## logcompactness_mean           1.000      0.885
## sqrtconcavity_mean            0.885      1.000
## sqrtconcave_points_mean      0.840      0.935
##                               sqrtconcave_points_mean
## radius_mean                   0.798
## texture_mean                  0.268
## perimeter_mean                 0.824
## smoothness_mean                0.563
## symmetry_mean                  0.438
## fractal_dimension_mean         0.158
## logarea_mean                   0.787
## logcompactness_mean             0.840
## sqrtconcavity_mean              0.935
## sqrtconcave_points_mean        1.000

library(corrplot)

## corrplot 0.92 loaded
wisconsin.corr = cor(wisconsin[, c(2:4, 6, 10:15)], method = c("spearman"))
corrplot(wisconsin.corr, order ='hclust', method="ellipse", addrect = 4)

```



As expected, mean log area, mean perimeter, and mean radius are highly correlated given that they are functions of each other ( $r > 0.982$ ). Many of our other predictors are correlated with each other. We will highlight some of the combinations that produce the strongest correlations. mean sqrt concavity and mean sqrt concave points are highly correlated with each other ( $r = 0.935$ ). This suggests that the severity of concave portions of the cell is correlated with the number of concave portions on the cell contour. Log transformed mean compactness is correlated with mean sqrt concavity and mean sqrt concave points ( $r = .885$  and  $.840$ ). As log transformed compactness increases so does mean sqrt concave points and mean sqrt concave points. There does not appear to be a logical explanation for why these predictors are correlated with each other. Issues with collinearity between the aforementioned correlated predictors should be considered when interpreting our output.

## First Order Logistic Regression Model

We will now fit a first order logistic regression model that includes all of the variables. Log transformed mean area and compactness were included instead of the original variables. Similarly, categorized mean concavity and concave points were included instead of the original variables.

```
lm1 <- glm(diagnosis ~ radius_mean + texture_mean + perimeter_mean + smoothness_mean + symmetry_mean + ...)
```

```

## 
## Call:
## glm(formula = diagnosis ~ radius_mean + texture_mean + perimeter_mean +
##      smoothness_mean + symmetry_mean + fractal_dimension_mean +
##      sqrtconcave_points_mean + sqrtconcavity_mean + logcompactness_mean +
##      logarea_mean, family = binomial, data = wisconsin)
##
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max
## -1.7933 -0.0980 -0.0124  0.0382  3.1847
##
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)           -146.54110   51.61835 -2.839  0.00453 **
## radius_mean            1.27412   3.09951  0.411  0.68102
## texture_mean           0.39790   0.06756  5.889 3.88e-09 ***
## perimeter_mean         -0.47562   0.52334 -0.909  0.36345
## smoothness_mean        89.97401  37.52975  2.397  0.01651 *
## symmetry_mean          21.79649  11.65929  1.869  0.06156 .
## fractal_dimension_mean 27.18567  78.64810  0.346  0.72960
## sqrtconcave_points_mean 18.78152  13.57631  1.383  0.16654
## sqrtconcavity_mean     18.38916  6.71775  2.737  0.00619 **
## logcompactness_mean     -2.28134   2.08386 -1.095  0.27362
## logarea_mean            21.05942  10.88736  1.934  0.05308 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 751.44 on 568 degrees of freedom
## Residual deviance: 140.66 on 558 degrees of freedom
## AIC: 162.66
##
## Number of Fisher Scoring iterations: 8
exp(0.39790)

## [1] 1.488695
exp(89.97401/100)

## [1] 2.458964
exp(21.79649/10)

## [1] 8.843202
exp(18.38916/10)

## [1] 6.289717
exp(21.05942/10)

## [1] 8.214838

```

A preliminary analysis of our coefficients based on our first order model:

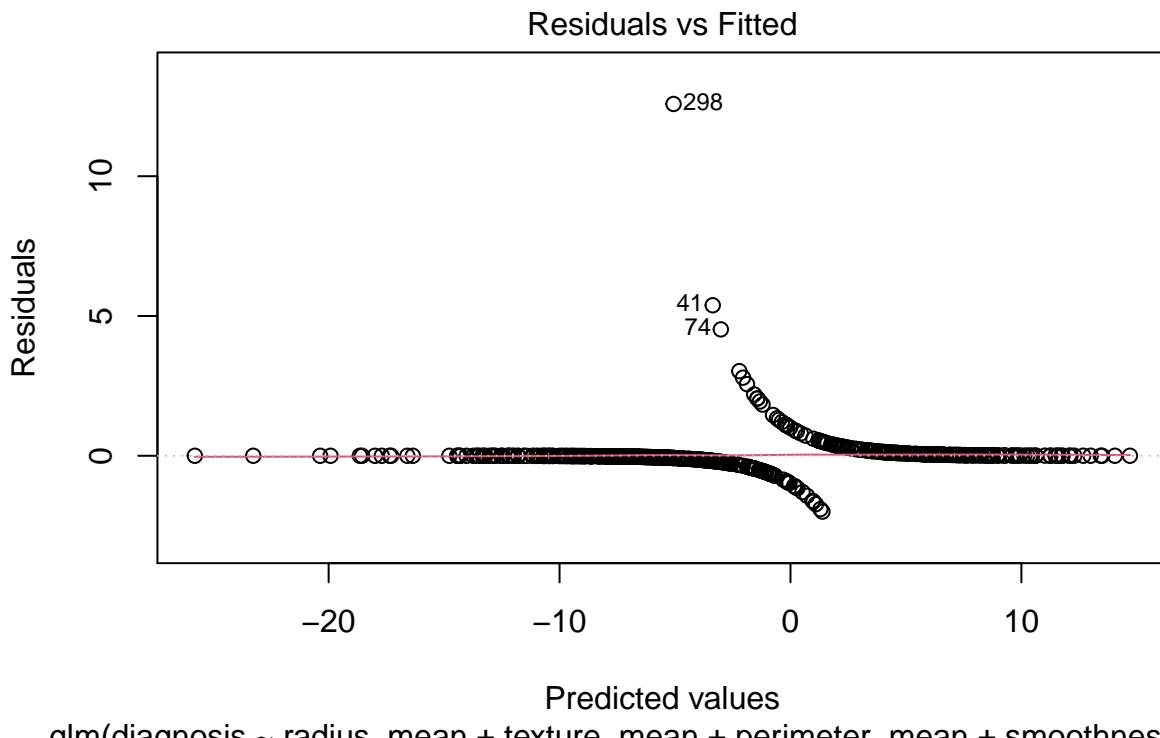
- Mean radius, perimeter, fractal dimension, sqrt concave points, and log compactness are not significantly associated with the odds of a tumor being malignant ( $p = 0.681, 0.363, 0.730, 0.167$ , and  $0.274$ )

respectively).

- For a one unit increase in mean cell texture, the odds of a tumor being malignant increases by a factor of 1.49, holding other predictors constant ( $p = 3.88e-09$ ).
- For an increase of 0.01 in mean cell smoothness, the odds of a tumor being malignant increases by a factor of 2.46, holding other predictors constant ( $p = 0.01651$ ).
- Holding other predictors constant, an increase of 0.1 in mean symmetry results in an increase in the odds of a tumor being malignant by a factor of 8.84, however this association is not quite significant ( $p = 0.06156$ ).
- For an increase of 0.1 in mean concavity, the odds of a tumor being malignant increases by a factor of 6.29, holding other predictors constant ( $p = 0.00619$ ).
- Holding other predictors constant, an increase of 0.1 in mean log area results in an increase in the odds of a tumor being malignant by a factor of 8.21, however this association is not quite significant ( $p = 0.05308$ ).

As previously mentioned, mean perimeter, radius, and log area are highly correlated. As a result, some of these correlated variables could be explaining similar variation in the log odds of a tumor being malignant. This should be taken into account when interpreting the coefficients and p values of these variables in the first order model.

```
par(mfrow=c(1,1))
plot(lm1, which = c(1))
```



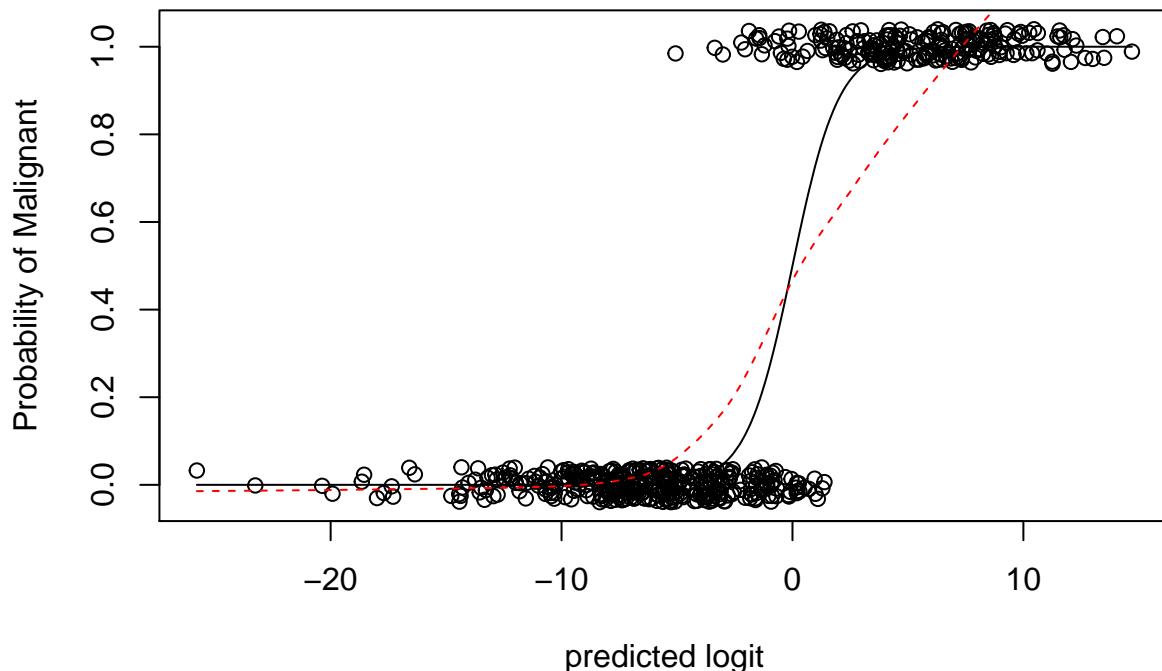
```
glm(diagnosis ~ radius_mean + texture_mean + perimeter_mean + smoothness_me
```

The smoothing spline in the residual plot for our first order model has little deviation from zero. Consequently, it appears that our model has successfully met the assumption of linearity.

```

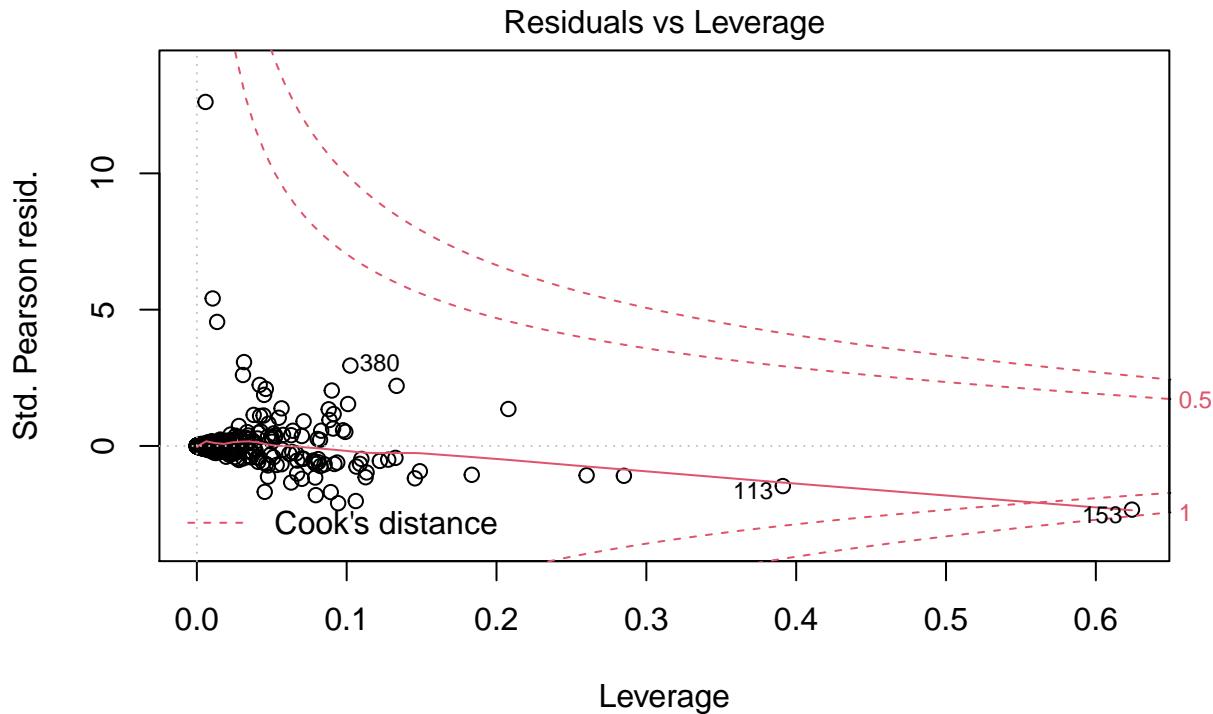
predpr = predict (lm1, type='response')
predlogit = predict (lm1)
plot (jitter (ifelse(wisconsin$diagnosis=='M',1,0), 0.2) ~ predlogit, xlab="predicted logit", ylab="Probability of Malignant")
pred.ord = order (predlogit)
lines (predlogit[pred.ord], predpr[pred.ord])
lines (lowess (predlogit [pred.ord], predpr [pred.ord]), col='red', lty=2)

```



By observing the density of points, there is a clear switch from benign to malignant tumor cells near the predicted logit of 0. The clear switch with little overlap from benign to malignant tumor cell responses as predicted logit increases suggests that our model is likely well suited to accurately predict tumor diagnosis. A fitted smoothing spline line fits the general pattern of the fitted logistic curve however it diverges slightly. This suggests that the logistic model may not be representing the trend in the data as well as it could.

```
plot(lm1, which=5)
```



A plot of residuals vs leverage on the first order model shows that case 153 and 113 have high leverage. Point 153 is also a moderate influence point. Point 298, while not influential due to low leverage, has a very high residual.

```
car::vif(lm1)
```

```
##          radius_mean      texture_mean      perimeter_mean
##      571.049501      1.784590      692.410572
##      smoothness_mean    symmetry_mean   fractal_dimension_mean
##      4.649794       1.927917       6.678347
##      sqrtconcave_points_mean    sqrtconcavity_mean   logcompactness_mean
##      7.039467        6.584806       11.814767
##      logarea_mean
##      141.186566
```

As could be expected from highly correlated predictors, our VIF values indicate that we have issues with collinearity. Particularly, we have collinearity issues with mean perimeter, radius, and log area (VIF values = 692.41, 571.05, and 141.19 respectively).

## Stepwise regression on the first-order model

```
lm2 = step(lm1, direction = 'both')

## Start:  AIC=162.66
## diagnosis ~ radius_mean + texture_mean + perimeter_mean + smoothness_mean +
##      symmetry_mean + fractal_dimension_mean + sqrtconcave_points_mean +
##      sqrtconcavity_mean + logcompactness_mean + logarea_mean
```

```

##                                     Df Deviance   AIC
## - fractal_dimension_mean      1  140.78 160.78
## - radius_mean                 1  140.83 160.83
## - perimeter_mean              1  141.46 161.46
## - logcompactness_mean         1  141.88 161.88
## - sqrtconcave_points_mean    1  142.65 162.65
## <none>                      140.66 162.66
## - symmetry_mean               1  144.19 164.19
## - logarea_mean                1  144.33 164.33
## - smoothness_mean              1  147.30 167.30
## - sqrtconcavity_mean          1  148.79 168.79
## - texture_mean                 1  190.05 210.05
##
## Step:  AIC=160.78
## diagnosis ~ radius_mean + texture_mean + perimeter_mean + smoothness_mean +
##             symmetry_mean + sqrtconcave_points_mean + sqrtconcavity_mean +
##             logcompactness_mean + logarea_mean
##
##                                     Df Deviance   AIC
## - radius_mean                  1  140.95 158.95
## - perimeter_mean               1  141.55 159.55
## - logcompactness_mean          1  141.89 159.89
## - sqrtconcave_points_mean     1  142.66 160.66
## <none>                      140.78 160.78
## - logarea_mean                 1  144.34 162.34
## - symmetry_mean                1  144.42 162.42
## + fractal_dimension_mean       1  140.66 162.66
## - smoothness_mean              1  148.07 166.07
## - sqrtconcavity_mean          1  150.37 168.37
## - texture_mean                 1  190.07 208.07
##
## Step:  AIC=158.95
## diagnosis ~ texture_mean + perimeter_mean + smoothness_mean +
##             symmetry_mean + sqrtconcave_points_mean + sqrtconcavity_mean +
##             logcompactness_mean + logarea_mean
##
##                                     Df Deviance   AIC
## - perimeter_mean               1  142.28 158.28
## - sqrtconcave_points_mean     1  142.75 158.75
## <none>                      140.95 158.95
## - logarea_mean                 1  144.48 160.48
## - symmetry_mean                1  144.62 160.62
## + radius_mean                  1  140.78 160.78
## + fractal_dimension_mean       1  140.83 160.83
## - logcompactness_mean          1  145.11 161.11
## - smoothness_mean              1  149.06 165.06
## - sqrtconcavity_mean          1  150.45 166.45
## - texture_mean                 1  190.43 206.43
##
## Step:  AIC=158.28
## diagnosis ~ texture_mean + smoothness_mean + symmetry_mean +
##             sqrtconcave_points_mean + sqrtconcavity_mean + logcompactness_mean +
##             logarea_mean

```

```

##                                     Df Deviance    AIC
## - sqrtconcave_points_mean   1   144.04 158.04
## <none>                      142.28 158.28
## + perimeter_mean            1   140.95 158.95
## + radius_mean               1   141.55 159.55
## - symmetry_mean             1   145.59 159.59
## + fractal_dimension_mean    1   142.25 160.25
## - sqrtconcavity_mean        1   151.21 165.21
## - smoothness_mean           1   151.42 165.42
## - logcompactness_mean       1   151.75 165.75
## - logarea_mean              1   181.26 195.26
## - texture_mean               1   193.05 207.05
##
## Step:  AIC=158.04
## diagnosis ~ texture_mean + smoothness_mean + symmetry_mean +
##           sqrtconcavity_mean + logcompactness_mean + logarea_mean
##
##                                     Df Deviance    AIC
## <none>                      144.04 158.04
## + sqrtconcave_points_mean   1   142.28 158.28
## + perimeter_mean            1   142.75 158.75
## - symmetry_mean             1   146.90 158.90
## + radius_mean               1   143.22 159.22
## + fractal_dimension_mean    1   144.04 160.04
## - logcompactness_mean       1   152.76 164.76
## - smoothness_mean           1   172.42 184.42
## - sqrtconcavity_mean        1   172.83 184.83
## - texture_mean               1   194.81 206.81
## - logarea_mean              1   302.62 314.62
summary(lm2)

##
## Call:
## glm(formula = diagnosis ~ texture_mean + smoothness_mean + symmetry_mean +
##       sqrtconcavity_mean + logcompactness_mean + logarea_mean,
##       family = binomial, data = wisconsin)
##
## Deviance Residuals:
##      Min        1Q     Median        3Q       Max
## -1.8187  -0.1161  -0.0200   0.0226   3.2896
##
## Coefficients:
##                                     Estimate Std. Error z value Pr(>|z|)
## (Intercept)                 -97.8137    12.7062 -7.698 1.38e-14 ***
## texture_mean                  0.3962     0.0658   6.022 1.73e-09 ***
## smoothness_mean                129.7312    27.5952   4.701 2.59e-06 ***
## symmetry_mean                  18.8677    11.2616   1.675  0.09385 .
## sqrtconcavity_mean              23.0600     4.5989   5.014 5.32e-07 ***
## logcompactness_mean             -3.2733     1.1556  -2.833  0.00462 **
## logarea_mean                   9.2641     1.2357   7.497 6.54e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##

```

```

## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 751.44 on 568 degrees of freedom
## Residual deviance: 144.04 on 562 degrees of freedom
## AIC: 158.04
##
## Number of Fisher Scoring iterations: 8

```

Our step wise regression using AIC criteria removed in order mean fractal dimension, radius, perimeter, and sqrt concave points. This resulted in a model that included mean texture, smoothness, sqrt concavity, log compactness, log area, and symmetry. All predictors were significant with the exception of mean symmetry for which the p value did not quite meet the 0.05 threshold ( $p = 0.094$ ). The AIC was minimized with this combination of predictors included in the model at 158.04. For comparison, the AIC of the full first order model was 162.66.

## Creating centered interaction effects

We will create centered interaction effects to hopefully counteract the previously identified collinearity issues.

```

my.center = function (x) (x - mean(x))
wisconsin$texture_mean.c = my.center(wisconsin$texture_mean)
wisconsin$smoothness_mean.c = my.center(wisconsin$smoothness_mean)
wisconsin$symmetry_mean.c = my.center(wisconsin$symmetry_mean)
wisconsin$sqrtconcavity_mean.c = my.center(wisconsin$sqrtconcavity_mean)
wisconsin$sqrtconcave_points_mean.c = my.center(wisconsin$concave points_mean)
wisconsin$logcompactness_mean.c = my.center(wisconsin$logcompactness_mean)
wisconsin$logarea_mean.c = my.center(wisconsin$logarea_mean)

```

## Fit a model with interactions

Only the predictors kept from the step wise regression on the first order model were centered and used in the interaction model to avoid convergence errors.

```

lm3 = glm(diagnosis ~ (texture_mean.c + smoothness_mean.c + symmetry_mean.c + sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c + logarea_mean.c)^2, family = binomial,
           data = wisconsin)

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
summary(lm3)

##
## Call:
## glm(formula = diagnosis ~ (texture_mean.c + smoothness_mean.c +
##     symmetry_mean.c + sqrtconcavity_mean.c + sqrtconcave_points_mean.c +
##     logcompactness_mean.c + logarea_mean.c)^2, family = binomial,
##     data = wisconsin)
##
## Deviance Residuals:
##      Min        1Q    Median        3Q       Max
## -1.9841   -0.0650   -0.0084    0.0002    3.4694
##
## Coefficients:
## (Intercept)          Estimate Std. Error z value
## texture_mean.c        4.878e-03  6.429e-01  0.008
## smoothness_mean.c     5.760e-01  1.697e-01  3.395
## symmetry_mean.c      1.219e+02  7.321e+01  1.666

```

```

## symmetry_mean.c          2.466e+01  3.043e+01  0.810
## sqrtconcavity_mean.c    4.909e+01  1.444e+01  3.400
## sqrtconcave_points_mean.c 6.181e+01  5.666e+01  1.091
## logcompactness_mean.c   -8.091e+00  2.885e+00 -2.805
## logarea_mean.c           9.067e+00  2.787e+00  3.254
## texture_mean.c:smoothness_mean.c 5.493e+01  2.237e+01  2.455
## texture_mean.c:symmetry_mean.c   4.995e+00  7.343e+00  0.680
## texture_mean.c:sqrtconcavity_mean.c 5.449e+00  3.559e+00  1.531
## texture_mean.c:sqrtconcave_points_mean.c -1.394e+01  1.692e+01 -0.824
## texture_mean.c:logcompactness_mean.c -1.482e+00  7.645e-01 -1.939
## texture_mean.c:logarea_mean.c     1.733e+00  6.933e-01  2.500
## smoothness_mean.c:symmetry_mean.c -6.726e+02  2.389e+03 -0.282
## smoothness_mean.c:sqrtconcavity_mean.c 5.193e+03  2.008e+03  2.586
## smoothness_mean.c:sqrtconcave_points_mean.c -1.205e+04  6.484e+03 -1.858
## smoothness_mean.c:logcompactness_mean.c -2.126e+02  3.201e+02 -0.664
## smoothness_mean.c:logarea_mean.c     -1.655e+02  1.529e+02 -1.083
## symmetry_mean.c:sqrtconcavity_mean.c -3.759e+02  7.131e+02 -0.527
## symmetry_mean.c:sqrtconcave_points_mean.c 2.283e+03  2.775e+03  0.823
## symmetry_mean.c:logcompactness_mean.c 5.014e+01  1.400e+02  0.358
## symmetry_mean.c:logarea_mean.c      -1.698e+01  9.548e+01 -0.178
## sqrtconcavity_mean.c:sqrtconcave_points_mean.c 3.321e+02  1.031e+03  0.322
## sqrtconcavity_mean.c:logcompactness_mean.c -6.203e+01  4.598e+01 -1.349
## sqrtconcavity_mean.c:logarea_mean.c    9.309e+01  5.842e+01  1.593
## sqrtconcave_points_mean.c:logcompactness_mean.c 7.175e+01  2.791e+02  0.257
## sqrtconcave_points_mean.c:logarea_mean.c -2.941e+02  2.264e+02 -1.299
## logcompactness_mean.c:logarea_mean.c    2.835e+00  1.073e+01  0.264
## Pr(>|z|)
## (Intercept)          0.993946
## texture_mean.c        0.000686 ***
## smoothness_mean.c    0.095807 .
## symmetry_mean.c      0.417783
## sqrtconcavity_mean.c 0.000675 ***
## sqrtconcave_points_mean.c 0.275292
## logcompactness_mean.c 0.005038 **
## logarea_mean.c        0.001139 **
## texture_mean.c:smoothness_mean.c 0.014072 *
## texture_mean.c:symmetry_mean.c   0.496401
## texture_mean.c:sqrtconcavity_mean.c 0.125707
## texture_mean.c:sqrtconcave_points_mean.c 0.410203
## texture_mean.c:logcompactness_mean.c 0.052512 .
## texture_mean.c:logarea_mean.c    0.012416 *
## smoothness_mean.c:symmetry_mean.c 0.778258
## smoothness_mean.c:sqrtconcavity_mean.c 0.009710 **
## smoothness_mean.c:sqrtconcave_points_mean.c 0.063159 .
## smoothness_mean.c:logcompactness_mean.c 0.506480
## smoothness_mean.c:logarea_mean.c    0.278941
## symmetry_mean.c:sqrtconcavity_mean.c 0.598137
## symmetry_mean.c:sqrtconcave_points_mean.c 0.410543
## symmetry_mean.c:logcompactness_mean.c 0.720131
## symmetry_mean.c:logarea_mean.c     0.858880
## sqrtconcavity_mean.c:sqrtconcave_points_mean.c 0.747365
## sqrtconcavity_mean.c:logcompactness_mean.c 0.177269
## sqrtconcavity_mean.c:logarea_mean.c   0.111098
## sqrtconcave_points_mean.c:logcompactness_mean.c 0.797088

```

```
## sqrtconcave_points_mean.c:logarea_mean.c          0.194012
## logcompactness_mean.c:logarea_mean.c            0.791645
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
## Null deviance: 751.440  on 568  degrees of freedom
## Residual deviance: 81.934  on 540  degrees of freedom
## AIC: 139.93
## 
## Number of Fisher Scoring iterations: 12
```

In this full second order model with centered predictors and interactions we find that mean texture, sqrt concavity, log compactness, and log compactness are significant predictors of tumor diagnosis. There are significant interactions between smoothness and texture, texture and log area, and between smoothness and sqrt concavity. The AIC of this model at 139.93 than the model generated from step wise regression of the full first order model without centered predictors at 158.04.

### Applying stepwise regression to the full second order model

AIC





```

## - texture_mean.c:symmetry_mean.c           1  82.541 136.54
## - texture_mean.c:sqrtconcave_points_mean.c 1  82.619 136.62
## - symmetry_mean.c:sqrtconcave_points_mean.c 1  82.625 136.62
## - smoothness_mean.c:logarea_mean.c          1  83.011 137.01
## - sqrtconcave_points_mean.c:logarea_mean.c   1  83.908 137.91
## <none>                                      1  81.963 137.96
## - sqrtconcavity_mean.c:logcompactness_mean.c 1  84.495 138.50
## - texture_mean.c:sqrtconcavity_mean.c        1  84.509 138.51
## - sqrtconcavity_mean.c:logarea_mean.c        1  85.350 139.35
## + symmetry_mean.c:logarea_mean.c            1  81.934 139.93
## - texture_mean.c:logcompactness_mean.c        1  87.018 141.02
## - smoothness_mean.c:sqrtconcave_points_mean.c 1  87.420 141.42
## - texture_mean.c:logarea_mean.c              1  89.489 143.49
## - texture_mean.c:smoothness_mean.c           1  91.018 145.02
## - smoothness_mean.c:sqrtconcavity_mean.c      1  98.274 152.27

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=136.01
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##           sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##           logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:symmetry_mean.c +
##           texture_mean.c:sqrtconcavity_mean.c + texture_mean.c:sqrtconcave_points_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:symmetry_mean.c + smoothness_mean.c:sqrtconcavity_mean.c +
##           smoothness_mean.c:sqrtconcave_points_mean.c + smoothness_mean.c:logcompactness_mean.c +
##           smoothness_mean.c:logarea_mean.c + symmetry_mean.c:sqrtconcavity_mean.c +
##           symmetry_mean.c:sqrtconcave_points_mean.c + symmetry_mean.c:logcompactness_mean.c +
##           sqrtconcavity_mean.c:sqrtconcave_points_mean.c + sqrtconcavity_mean.c:logcompactness_mean.c +
##           sqrtconcavity_mean.c:logarea_mean.c + sqrtconcave_points_mean.c:logcompactness_mean.c +
##           sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## - sqrtconcave_points_mean.c:logcompactness_mean.c 1 82.296 132.30
## - texture_mean.c:symmetry_mean.c 1 82.598 132.60
## - texture_mean.c:sqrtconcave_points_mean.c 1 82.682 132.68
## - symmetry_mean.c:sqrtconcave_points_mean.c 1 82.690 132.69
## - smoothness_mean.c:logarea_mean.c 1 83.091 133.09
## - smoothness_mean.c:logcompactness_mean.c 1 83.118 133.12
## <none> 82.057 134.06
## - sqrtconcave_points_mean.c:logarea_mean.c 1 84.112 134.11
## - texture_mean.c:sqrtconcavity_mean.c 1 84.613 134.61
## - sqrtconcavity_mean.c:logcompactness_mean.c 1 84.861 134.86
## + smoothness_mean.c:symmetry_mean.c 1 82.013 136.01
## + logcompactness_mean.c:logarea_mean.c 1 82.015 136.01
## + symmetry_mean.c:logarea_mean.c 1 82.057 136.06
## - texture_mean.c:logcompactness_mean.c 1 87.319 137.32
## - smoothness_mean.c:sqrtconcave_points_mean.c 1 88.173 138.17
## - sqrtconcavity_mean.c:logarea_mean.c 1 88.464 138.46
## - texture_mean.c:logarea_mean.c 1 89.685 139.69
## - texture_mean.c:smoothness_mean.c 1 91.176 141.18
## - smoothness_mean.c:sqrtconcavity_mean.c 1 100.635 150.63

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=132.09
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##           sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##           logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:symmetry_mean.c +
##           texture_mean.c:sqrtconcavity_mean.c + texture_mean.c:sqrtconcave_points_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##           smoothness_mean.c:logcompactness_mean.c + smoothness_mean.c:logarea_mean.c +
##           symmetry_mean.c:sqrtconcavity_mean.c + symmetry_mean.c:sqrtconcave_points_mean.c +
##           symmetry_mean.c:logcompactness_mean.c + sqrtconcavity_mean.c:logcompactness_mean.c +
##           sqrtconcavity_mean.c:logarea_mean.c + sqrtconcave_points_mean.c:logcompactness_mean.c +
##           sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## - sqrtconcave_points_mean.c:logcompactness_mean.c 1 82.611 128.61
## - texture_mean.c:sqrtconcave_points_mean.c 1 82.784 128.78
## - texture_mean.c:symmetry_mean.c 1 82.797 128.80
## - symmetry_mean.c:sqrtconcave_points_mean.c 1 83.007 129.01
## - smoothness_mean.c:logarea_mean.c 1 83.233 129.23
## <none> 82.122 130.12
## - sqrtconcave_points_mean.c:logarea_mean.c 1 84.297 130.30
## - smoothness_mean.c:logcompactness_mean.c 1 84.402 130.40
## - sqrtconcavity_mean.c:logcompactness_mean.c 1 84.893 130.89
## - texture_mean.c:sqrtconcavity_mean.c 1 85.032 131.03
## + symmetry_mean.c:logcompactness_mean.c 1 82.090 132.09
## + sqrtconcavity_mean.c:sqrtconcave_points_mean.c 1 82.110 132.11
## + symmetry_mean.c:logarea_mean.c 1 82.113 132.11
## + logcompactness_mean.c:logarea_mean.c 1 82.114 132.11
## + smoothness_mean.c:symmetry_mean.c 1 82.119 132.12
## - texture_mean.c:logcompactness_mean.c 1 88.438 134.44
## - smoothness_mean.c:sqrtconcave_points_mean.c 1 89.510 135.51
## - texture_mean.c:logarea_mean.c 1 89.807 135.81
## - sqrtconcavity_mean.c:logarea_mean.c 1 90.698 136.70
## - texture_mean.c:smoothness_mean.c 1 91.629 137.63
## - smoothness_mean.c:sqrtconcavity_mean.c 1 102.310 148.31

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=128.24
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##     sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##     logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:symmetry_mean.c +
##     texture_mean.c:sqrtconcavity_mean.c + texture_mean.c:sqrtconcave_points_mean.c +
##     texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##     smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##     smoothness_mean.c:logcompactness_mean.c + smoothness_mean.c:logarea_mean.c +
##     symmetry_mean.c:sqrtconcave_points_mean.c + sqrtconcavity_mean.c:logcompactness_mean.c +
##     sqrtconcavity_mean.c:logarea_mean.c + sqrtconcave_points_mean.c:logcompactness_mean.c +
##     sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## - sqrtconcave_points_mean.c:logcompactness_mean.c 1 83.895 125.89
## - symmetry_mean.c:sqrtconcave_points_mean.c      1 84.195 126.19
## - smoothness_mean.c:logarea_mean.c               1 84.456 126.46
## <none>
## - smoothness_mean.c:logcompactness_mean.c        1 85.327 127.33
## + texture_mean.c:symmetry_mean.c                 1 82.237 128.24
## - sqrtconcave_points_mean.c:logarea_mean.c       1 86.472 128.47
## + symmetry_mean.c:sqrtconcavity_mean.c          1 82.797 128.80
## + symmetry_mean.c:logarea_mean.c                 1 82.816 128.82
## + symmetry_mean.c:logcompactness_mean.c          1 82.829 128.83
## + sqrtconcavity_mean.c:sqrtconcave_points_mean.c 1 82.853 128.85
## + smoothness_mean.c:symmetry_mean.c              1 82.860 128.86
## + logcompactness_mean.c:logarea_mean.c           1 82.861 128.86
## - texture_mean.c:sqrtconcavity_mean.c            1 88.158 130.16
## - sqrtconcavity_mean.c:logcompactness_mean.c     1 88.650 130.65
## - texture_mean.c:logcompactness_mean.c            1 88.913 130.91
## - texture_mean.c:logarea_mean.c                  1 90.112 132.11
## - smoothness_mean.c:sqrtconcave_points_mean.c    1 91.215 133.22
## - texture_mean.c:smoothness_mean.c                1 94.619 136.62
## - sqrtconcavity_mean.c:logarea_mean.c            1 97.787 139.79
## - smoothness_mean.c:sqrtconcavity_mean.c          1 105.900 147.90

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=125.83
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##           sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##           logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:sqrtconcavity_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##           smoothness_mean.c:logcompactness_mean.c + smoothness_mean.c:logarea_mean.c +
##           symmetry_mean.c:sqrtconcave_points_mean.c + sqrtconcavity_mean.c:logcompactness_mean.c +
##           sqrtconcavity_mean.c:logarea_mean.c + sqrtconcave_points_mean.c:logcompactness_mean.c +
##           sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## - smoothness_mean.c:logcompactness_mean.c      1  86.671 124.67
## <none>                                         1  84.681 124.68
## + texture_mean.c:symmetry_mean.c             1  83.523 125.52
## + sqrtconcave_points_mean.c:logcompactness_mean.c 1  83.829 125.83
## + texture_mean.c:sqrtconcave_points_mean.c     1  83.895 125.89
## + symmetry_mean.c:sqrtconcavity_mean.c        1  84.333 126.33
## + sqrtconcavity_mean.c:sqrtconcave_points_mean.c 1  84.389 126.39
## + smoothness_mean.c:symmetry_mean.c            1  84.638 126.64
## + logcompactness_mean.c:logarea_mean.c         1  84.679 126.68
## + symmetry_mean.c:logarea_mean.c              1  84.680 126.68
## + symmetry_mean.c:logcompactness_mean.c        1  84.680 126.68
## - texture_mean.c:sqrtconcavity_mean.c          1  89.842 127.84
## - sqrtconcavity_mean.c:logcompactness_mean.c   1  90.346 128.35
## - smoothness_mean.c:sqrtconcave_points_mean.c  1  91.371 129.37
## - texture_mean.c:logarea_mean.c                1  93.086 131.09
## - texture_mean.c:logcompactness_mean.c          1  95.507 133.51
## - sqrtconcavity_mean.c:logarea_mean.c          1  98.398 136.40
## - texture_mean.c:smoothness_mean.c             1 103.921 141.92
## - smoothness_mean.c:sqrtconcavity_mean.c        1 106.847 144.85

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=124.39
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##           sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##           logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:sqrtconcavity_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##           smoothness_mean.c:logcompactness_mean.c + symmetry_mean.c:sqrtconcave_points_mean.c +
##           sqrtconcavity_mean.c:logcompactness_mean.c + sqrtconcavity_mean.c:logarea_mean.c +
##           sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## + smoothness_mean.c:logarea_mean.c           1  86.671 124.67
## + smoothness_mean.c:symmetry_mean.c         1  87.364 125.36
## + sqrtconcave_points_mean.c:logcompactness_mean.c 1  87.378 125.38
## + symmetry_mean.c:logcompactness_mean.c      1  87.467 125.47
## + symmetry_mean.c:sqrtconcavity_mean.c       1  87.518 125.52
## + logcompactness_mean.c:logarea_mean.c        1  87.541 125.54
## + symmetry_mean.c:logarea_mean.c             1  87.583 125.58
## - texture_mean.c:sqrtconcavity_mean.c        1  92.066 126.07
## - smoothness_mean.c:sqrtconcave_points_mean.c 1  95.391 129.39
## - sqrtconcavity_mean.c:logcompactness_mean.c 1  95.867 129.87
## - texture_mean.c:logarea_mean.c              1  96.860 130.86
## - texture_mean.c:logcompactness_mean.c        1  97.656 131.66
## - sqrtconcavity_mean.c:logarea_mean.c        1  99.340 133.34
## - texture_mean.c:smoothness_mean.c            1 106.059 140.06
## - smoothness_mean.c:sqrtconcavity_mean.c      1 109.846 143.85

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=123.33
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##           sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##           logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:sqrtconcavity_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##           sqrtconcavity_mean.c:logcompactness_mean.c + sqrtconcavity_mean.c:logarea_mean.c +
##           sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```







```

## - texture_mean.c:logarea_mean.c           1  96.251 128.25
## - sqrtconcavity_mean.c:logcompactness_mean.c 1  96.469 128.47
## - smoothness_mean.c:sqrtconcave_points_mean.c 1  96.552 128.55
## - texture_mean.c:logcompactness_mean.c       1 100.972 132.97
## - sqrtconcavity_mean.c:logarea_mean.c        1 100.985 132.99
## - texture_mean.c:smoothness_mean.c          1 107.055 139.06
## - smoothness_mean.c:sqrtconcavity_mean.c     1 111.371 143.37

summary(lm4)

##
## Call:
## glm(formula = diagnosis ~ texture_mean.c + smoothness_mean.c +
##      sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##      logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:sqrtconcavity_mean.c +
##      texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##      smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##      sqrtconcavity_mean.c:logcompactness_mean.c + sqrtconcavity_mean.c:logarea_mean.c +
##      sqrtconcave_points_mean.c:logarea_mean.c + sqrtconcavity_mean.c:sqrtconcave_points_mean.c,
##      family = binomial, data = wisconsin)
##
## Deviance Residuals:
##    Min      1Q   Median      3Q      Max
## -2.3712 -0.0848 -0.0091  0.0002  3.6723
##
## Coefficients:
##                               Estimate Std. Error z value
## (Intercept)                3.028e-02 5.415e-01  0.056
## texture_mean.c              6.319e-01 1.235e-01  5.118
## smoothness_mean.c           1.288e+02 6.236e+01  2.066
## sqrtconcavity_mean.c        5.112e+01 1.276e+01  4.006
## sqrtconcave_points_mean.c  5.799e+01 4.297e+01  1.349
## logcompactness_mean.c       -7.684e+00 2.256e+00 -3.407
## logarea_mean.c              9.597e+00 2.392e+00  4.013
## texture_mean.c:smoothness_mean.c 4.827e+01 1.467e+01  3.290
## texture_mean.c:sqrtconcavity_mean.c 5.212e+00 2.128e+00  2.450
## texture_mean.c:logcompactness_mean.c -1.650e+00 5.182e-01 -3.184
## texture_mean.c:logarea_mean.c      1.307e+00 4.851e-01  2.693
## smoothness_mean.c:sqrtconcavity_mean.c 4.045e+03 1.176e+03  3.441
## smoothness_mean.c:sqrtconcave_points_mean.c -9.822e+03 3.878e+03 -2.533
## sqrtconcavity_mean.c:logcompactness_mean.c -5.896e+01 2.473e+01 -2.385
## sqrtconcavity_mean.c:logarea_mean.c      9.540e+01 3.055e+01  3.123
## sqrtconcave_points_mean.c:logarea_mean.c -3.080e+02 1.388e+02 -2.220
## sqrtconcavity_mean.c:sqrtconcave_points_mean.c 6.967e+02 4.428e+02  1.573
##                               Pr(>|z|)
## (Intercept)                0.955407
## texture_mean.c              3.08e-07 ***
## smoothness_mean.c            0.038870 *
## sqrtconcavity_mean.c         6.17e-05 ***
## sqrtconcave_points_mean.c   0.177194
## logcompactness_mean.c        0.000658 ***
## logarea_mean.c               6.00e-05 ***
## texture_mean.c:smoothness_mean.c 0.001003 **
## texture_mean.c:sqrtconcavity_mean.c 0.014305 *
## texture_mean.c:logcompactness_mean.c 0.001454 **
```

```

## texture_mean.c:logarea_mean.c          0.007071 **
## smoothness_mean.c:sqrtconcavity_mean.c 0.000579 ***
## smoothness_mean.c:sqrtconcave_points_mean.c 0.011319 *
## sqrtconcavity_mean.c:logcompactness_mean.c 0.017100 *
## sqrtconcavity_mean.c:logarea_mean.c      0.001791 **
## sqrtconcave_points_mean.c:logarea_mean.c 0.026442 *
## sqrtconcavity_mean.c:sqrtconcave_points_mean.c 0.115656
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 751.440  on 568  degrees of freedom
## Residual deviance: 87.729  on 552  degrees of freedom
## AIC: 121.73
##
## Number of Fisher Scoring iterations: 11

```

Given that our model generated using the AIC criteria kept a relatively large number of predictors (6) and interaction effects (10), we will run another step wise regression using BIC criteria which is more restrictive about which predictor variables are retained and compare the outputs

BIC

```

lm5 = step(lm3, direction = "both", k = log(lm3$rank + lm3$df.residual))

## Start:  AIC=265.91
## diagnosis ~ (texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##               sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##               logarea_mean.c)^2

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## - texture_mean.c:symmetry_mean.c           1  82.541 253.83
## - texture_mean.c:sqrtconcave_points_mean.c 1  82.619 253.90
## - symmetry_mean.c:sqrtconcave_points_mean.c 1  82.625 253.91
## - smoothness_mean.c:logarea_mean.c          1  83.011 254.30
## - sqrtconcave_points_mean.c:logarea_mean.c   1  83.908 255.19
## - sqrtconcavity_mean.c:logcompactness_mean.c 1  84.495 255.78
## - texture_mean.c:sqrtconcavity_mean.c        1  84.509 255.79
## - sqrtconcavity_mean.c:logarea_mean.c         1  85.350 256.63
## - texture_mean.c:logcompactness_mean.c        1  87.018 258.30
## - smoothness_mean.c:sqrtconcave_points_mean.c 1  87.420 258.70
## <none>                                         81.963 259.59
## - texture_mean.c:logarea_mean.c              1  89.489 260.77
## - texture_mean.c:smoothness_mean.c            1  91.018 262.30
## + symmetry_mean.c:logarea_mean.c             1  81.934 265.91
## - smoothness_mean.c:sqrtconcavity_mean.c      1  98.274 269.56

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=253.3
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##           sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##           logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:symmetry_mean.c +
##           texture_mean.c:sqrtconcavity_mean.c + texture_mean.c:sqrtconcave_points_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:symmetry_mean.c + smoothness_mean.c:sqrtconcavity_mean.c +
##           smoothness_mean.c:sqrtconcave_points_mean.c + smoothness_mean.c:logcompactness_mean.c +
##           smoothness_mean.c:logarea_mean.c + symmetry_mean.c:sqrtconcavity_mean.c +
##           symmetry_mean.c:sqrtconcave_points_mean.c + symmetry_mean.c:logcompactness_mean.c +
##           sqrtconcavity_mean.c:sqrtconcave_points_mean.c + sqrtconcavity_mean.c:logcompactness_mean.c +
##           sqrtconcavity_mean.c:logarea_mean.c + sqrtconcave_points_mean.c:logcompactness_mean.c +
##           sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## - sqrtconcave_points_mean.c:logcompactness_mean.c 1 82.296 240.89
## - texture_mean.c:symmetry_mean.c 1 82.598 241.19
## - texture_mean.c:sqrtconcave_points_mean.c 1 82.682 241.28
## - symmetry_mean.c:sqrtconcave_points_mean.c 1 82.690 241.29
## - smoothness_mean.c:logarea_mean.c 1 83.091 241.69
## - smoothness_mean.c:logcompactness_mean.c 1 83.118 241.72
## - sqrtconcave_points_mean.c:logarea_mean.c 1 84.112 242.71
## - texture_mean.c:sqrtconcavity_mean.c 1 84.613 243.21
## - sqrtconcavity_mean.c:logcompactness_mean.c 1 84.861 243.46
## - texture_mean.c:logcompactness_mean.c 1 87.319 245.92
## - smoothness_mean.c:sqrtconcave_points_mean.c 1 88.173 246.77
## <none> 82.057 247.00
## - sqrtconcavity_mean.c:logarea_mean.c 1 88.464 247.06
## - texture_mean.c:logarea_mean.c 1 89.685 248.28
## - texture_mean.c:smoothness_mean.c 1 91.176 249.77
## + smoothness_mean.c:symmetry_mean.c 1 82.013 253.30
## + logcompactness_mean.c:logarea_mean.c 1 82.015 253.30
## + symmetry_mean.c:logarea_mean.c 1 82.057 253.34
## - smoothness_mean.c:sqrtconcavity_mean.c 1 100.635 259.23

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=240.69
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##           sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##           logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:symmetry_mean.c +
##           texture_mean.c:sqrtconcavity_mean.c + texture_mean.c:sqrtconcave_points_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##           smoothness_mean.c:logcompactness_mean.c + smoothness_mean.c:logarea_mean.c +
##           symmetry_mean.c:sqrtconcavity_mean.c + symmetry_mean.c:sqrtconcave_points_mean.c +
##           symmetry_mean.c:logcompactness_mean.c + sqrtconcavity_mean.c:logcompactness_mean.c +
##           sqrtconcavity_mean.c:logarea_mean.c + sqrtconcave_points_mean.c:logcompactness_mean.c +
##           sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## - sqrtconcave_points_mean.c:logcompactness_mean.c 1 82.611 228.52
## - texture_mean.c:sqrtconcave_points_mean.c 1 82.784 228.69
## - texture_mean.c:symmetry_mean.c 1 82.797 228.71
## - symmetry_mean.c:sqrtconcave_points_mean.c 1 83.007 228.92
## - smoothness_mean.c:logarea_mean.c 1 83.233 229.14
## - sqrtconcave_points_mean.c:logarea_mean.c 1 84.297 230.21
## - smoothness_mean.c:logcompactness_mean.c 1 84.402 230.31
## - sqrtconcavity_mean.c:logcompactness_mean.c 1 84.893 230.80
## - texture_mean.c:sqrtconcavity_mean.c 1 85.032 230.94
## - texture_mean.c:logcompactness_mean.c 1 88.438 234.35
## <none> 82.122 234.38
## - smoothness_mean.c:sqrtconcave_points_mean.c 1 89.510 235.42
## - texture_mean.c:logarea_mean.c 1 89.807 235.72
## - sqrtconcavity_mean.c:logarea_mean.c 1 90.698 236.61
## - texture_mean.c:smoothness_mean.c 1 91.629 237.54
## + symmetry_mean.c:logcompactness_mean.c 1 82.090 240.69
## + sqrtconcavity_mean.c:sqrtconcave_points_mean.c 1 82.110 240.71
## + symmetry_mean.c:logarea_mean.c 1 82.113 240.71
## + logcompactness_mean.c:logarea_mean.c 1 82.114 240.71
## + smoothness_mean.c:symmetry_mean.c 1 82.119 240.72
## - smoothness_mean.c:sqrtconcavity_mean.c 1 102.310 248.22

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=228.15
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##     sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##     logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:symmetry_mean.c +
##     texture_mean.c:sqrtconcavity_mean.c + texture_mean.c:sqrtconcave_points_mean.c +
##     texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##     smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##     smoothness_mean.c:logcompactness_mean.c + smoothness_mean.c:logarea_mean.c +
##     symmetry_mean.c:sqrtconcave_points_mean.c + sqrtconcavity_mean.c:logcompactness_mean.c +
##     sqrtconcavity_mean.c:logarea_mean.c + sqrtconcave_points_mean.c:logcompactness_mean.c +
##     sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## - sqrtconcave_points_mean.c:logcompactness_mean.c 1 83.895 217.12
## - symmetry_mean.c:sqrtconcave_points_mean.c      1 84.195 217.42
## - smoothness_mean.c:logarea_mean.c               1 84.456 217.68
## - smoothness_mean.c:logcompactness_mean.c        1 85.327 218.55
## - sqrtconcave_points_mean.c:logarea_mean.c       1 86.472 219.69
## - texture_mean.c:sqrtconcavity_mean.c           1 88.158 221.38
## - sqrtconcavity_mean.c:logcompactness_mean.c     1 88.650 221.87
## - texture_mean.c:logcompactness_mean.c           1 88.913 222.13
## <none>                                         82.865 222.43
## - texture_mean.c:logarea_mean.c                 1 90.112 223.33
## - smoothness_mean.c:sqrtconcave_points_mean.c    1 91.215 224.44
## - texture_mean.c:smoothness_mean.c              1 94.619 227.84
## + texture_mean.c:symmetry_mean.c                1 82.237 228.15
## + symmetry_mean.c:sqrtconcavity_mean.c         1 82.797 228.71
## + symmetry_mean.c:logarea_mean.c                1 82.816 228.73
## + symmetry_mean.c:logcompactness_mean.c         1 82.829 228.74
## + sqrtconcavity_mean.c:sqrtconcave_points_mean.c 1 82.853 228.76
## + smoothness_mean.c:symmetry_mean.c             1 82.860 228.77
## + logcompactness_mean.c:logarea_mean.c          1 82.861 228.77
## - sqrtconcavity_mean.c:logarea_mean.c           1 97.787 231.01
## - smoothness_mean.c:sqrtconcavity_mean.c        1 105.900 239.12

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=217.05
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##           sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##           logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:sqrtconcavity_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##           smoothness_mean.c:logcompactness_mean.c + smoothness_mean.c:logarea_mean.c +
##           symmetry_mean.c:sqrtconcave_points_mean.c + sqrtconcavity_mean.c:logcompactness_mean.c +
##           sqrtconcavity_mean.c:logarea_mean.c + sqrtconcave_points_mean.c:logcompactness_mean.c +
##           sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## - smoothness_mean.c:logcompactness_mean.c      1  86.671 207.21
## - texture_mean.c:sqrtconcavity_mean.c         1  89.842 210.38
## - sqrtconcavity_mean.c:logcompactness_mean.c   1  90.346 210.88
## <none>                                         84.681 211.56
## - smoothness_mean.c:sqrtconcave_points_mean.c 1  91.371 211.91
## - texture_mean.c:logarea_mean.c                1  93.086 213.62
## - texture_mean.c:logcompactness_mean.c          1  95.507 216.04
## + texture_mean.c:symmetry_mean.c               1  83.523 216.75
## + sqrtconcave_points_mean.c:logcompactness_mean.c 1  83.829 217.05
## + texture_mean.c:sqrtconcave_points_mean.c     1  83.895 217.12
## + symmetry_mean.c:sqrtconcavity_mean.c         1  84.333 217.55
## + sqrtconcavity_mean.c:sqrtconcave_points_mean.c 1  84.389 217.61
## + smoothness_mean.c:symmetry_mean.c            1  84.638 217.86
## + logcompactness_mean.c:logarea_mean.c          1  84.679 217.90
## + symmetry_mean.c:logarea_mean.c               1  84.680 217.90
## + symmetry_mean.c:logcompactness_mean.c         1  84.680 217.90
## - sqrtconcavity_mean.c:logarea_mean.c          1  98.398 218.93
## - texture_mean.c:smoothness_mean.c              1 103.921 224.46
## - smoothness_mean.c:sqrtconcavity_mean.c        1 106.847 227.38

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=206.93
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##           sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##           logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:sqrtconcavity_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##           smoothness_mean.c:logcompactness_mean.c + symmetry_mean.c:sqrtconcave_points_mean.c +
##           sqrtconcavity_mean.c:logcompactness_mean.c + sqrtconcavity_mean.c:logarea_mean.c +
##           sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```





```

## - texture_mean.c:logcompactness_mean.c      1  97.656 205.50
## + texture_mean.c:symmetry_mean.c           1  85.965 206.50
## + sqrtconcavity_mean.c:sqrtconcave_points_mean.c 1  86.043 206.58
## + smoothness_mean.c:logcompactness_mean.c    1  86.393 206.93
## + texture_mean.c:sqrtconcave_points_mean.c   1  86.553 207.09
## - sqrtconcavity_mean.c:logarea_mean.c       1  99.340 207.19
## + smoothness_mean.c:logarea_mean.c          1  86.671 207.21
## + smoothness_mean.c:symmetry_mean.c         1  87.364 207.90
## + sqrtconcave_points_mean.c:logcompactness_mean.c 1  87.378 207.91
## + symmetry_mean.c:logcompactness_mean.c     1  87.467 208.00
## + symmetry_mean.c:sqrtconcavity_mean.c      1  87.518 208.05
## + logcompactness_mean.c:logarea_mean.c       1  87.541 208.07
## + symmetry_mean.c:logarea_mean.c            1  87.583 208.12
## - texture_mean.c:smoothness_mean.c          1  106.059 213.91
## - smoothness_mean.c:sqrtconcavity_mean.c     1  109.846 217.69

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=197.18
## diagnosis ~ texture_mean.c + smoothness_mean.c + symmetry_mean.c +
##           sqrtconcavity_mean.c + sqrtconcave_points_mean.c + logcompactness_mean.c +
##           logarea_mean.c + texture_mean.c:smoothness_mean.c + texture_mean.c:sqrtconcavity_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:sqrtconcavity_mean.c + smoothness_mean.c:sqrtconcave_points_mean.c +
##           sqrtconcavity_mean.c:logcompactness_mean.c + sqrtconcavity_mean.c:logarea_mean.c +
##           sqrtconcave_points_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```







```

## + texture_mean.c:sqrtconcave_points_mean.c      1  93.409 194.91
## + sqrtconcave_points_mean.c:logcompactness_mean.c 1  93.431 194.93
## - texture_mean.c:logcompactness_mean.c           1 106.564 195.38
## - texture_mean.c:logarea_mean.c                 1 109.065 197.88
## - smoothness_mean.c:sqrtconcavity_mean.c        1 111.763 200.58
## - texture_mean.c:smoothness_mean.c              1 114.985 203.80

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=186.35
## diagnosis ~ texture_mean.c + smoothness_mean.c + sqrtconcavity_mean.c +
##           sqrtconcave_points_mean.c + logcompactness_mean.c + logarea_mean.c +
##           texture_mean.c:smoothness_mean.c + texture_mean.c:sqrtconcavity_mean.c +
##           texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##           smoothness_mean.c:sqrtconcavity_mean.c + sqrtconcavity_mean.c:logcompactness_mean.c +
##           sqrtconcavity_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## - sqrtconcave_points_mean.c                      Df Deviance    AIC
##                                         1  98.024 180.50

```

```

## - sqrtconcavity_mean.c:logcompactness_mean.c      1 101.226 183.70
## - sqrtconcavity_mean.c:logarea_mean.c           1 103.435 185.91
## <none>                                         97.533 186.35
## - texture_mean.c:sqrtconcavity_mean.c          1 104.280 186.75
## + smoothness_mean.c:sqrtconcave_points_mean.c   1 93.490 188.65
## + smoothness_mean.c:logarea_mean.c             1 95.509 190.67
## + smoothness_mean.c:logcompactness_mean.c       1 95.520 190.68
## + symmetry_mean.c                            1 96.234 191.39
## + sqrtconcave_points_mean.c:logcompactness_mean.c 1 96.250 191.41
## + logcompactness_mean.c:logarea_mean.c          1 96.540 191.70
## + sqrtconcave_points_mean.c:logarea_mean.c       1 96.627 191.78
## + texture_mean.c:sqrtconcave_points_mean.c     1 97.531 192.69
## + sqrtconcavity_mean.c:sqrtconcave_points_mean.c 1 97.531 192.69
## - texture_mean.c:logcompactness_mean.c          1 111.611 194.08
## - texture_mean.c:logarea_mean.c                1 113.852 196.32
## - smoothness_mean.c:sqrtconcavity_mean.c        1 118.693 201.16
## - texture_mean.c:smoothness_mean.c              1 129.317 211.79

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##
## Step: AIC=180.49
## diagnosis ~ texture_mean.c + smoothness_mean.c + sqrtconcavity_mean.c +
##           logcompactness_mean.c + logarea_mean.c + texture_mean.c:smoothness_mean.c +
##           texture_mean.c:sqrtconcavity_mean.c + texture_mean.c:logcompactness_mean.c +
##           texture_mean.c:logarea_mean.c + smoothness_mean.c:sqrtconcavity_mean.c +
##           sqrtconcavity_mean.c:logcompactness_mean.c + sqrtconcavity_mean.c:logarea_mean.c

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## Df Deviance    AIC
## - sqrtconcavity_mean.c:logcompactness_mean.c 1 101.384 177.51

```





```

## - texture_mean.c:smoothness_mean.c           1   132.65 196.09
summary(lm5)

##
## Call:
## glm(formula = diagnosis ~ texture_mean.c + smoothness_mean.c +
##      sqrtconcavity_mean.c + logcompactness_mean.c + logarea_mean.c +
##      texture_mean.c:smoothness_mean.c + texture_mean.c:sqrtconcavity_mean.c +
##      texture_mean.c:logcompactness_mean.c + texture_mean.c:logarea_mean.c +
##      smoothness_mean.c:sqrtconcavity_mean.c, family = binomial,
##      data = wisconsin)
##
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max
## -2.2121 -0.0999 -0.0137  0.0009  3.5626
##
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)                -1.2176    0.3238 -3.760 0.000170 ***
## texture_mean.c               0.5738    0.1047  5.481 4.23e-08 ***
## smoothness_mean.c            206.4960   43.8422  4.710 2.48e-06 ***
## sqrtconcavity_mean.c        29.9367   6.2486  4.791 1.66e-06 ***
## logcompactness_mean.c       -4.1492   1.5654 -2.651 0.008035 **
## logarea_mean.c              12.2830   1.8590  6.607 3.91e-11 ***
## texture_mean.c:smoothness_mean.c 41.2347   9.8465  4.188 2.82e-05 ***
## texture_mean.c:sqrtconcavity_mean.c 4.7171   1.4489  3.256 0.001131 **
## texture_mean.c:logcompactness_mean.c -1.3355   0.3649 -3.660 0.000252 ***
## texture_mean.c:logarea_mean.c      1.0290   0.3774  2.727 0.006397 **
## smoothness_mean.c:sqrtconcavity_mean.c 967.4881  276.3854  3.501 0.000464 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 751.44 on 568 degrees of freedom
## Residual deviance: 107.18 on 558 degrees of freedom
## AIC: 129.18
##
## Number of Fisher Scoring iterations: 9

```

The BIC model kept mean texture, smoothness, sqrt concavity, log compactness, log area, and the interaction between texture and smoothness, sqrt concavity, log compactness, and log area as well as the interaction between smoothness and sqrt concavity. Overall, 5 predictors and 5 interactions, all of which were significant, were kept in the final model. Residual deviance (107.18 in BIC model and 87.729 AIC model) and AIC values (129.18 in BIC model and 121.73 in AIC model) are relatively similar between the BIC and the AIC models. This gives us confidence proceeding with the BIC model as our final model for simplicity's sake. The AIC of this model (129.18) is considerably smaller than the AIC of the full first order model (162.66) or the model generated from the step wise regression on the first order model (158.04).

## Interaction plots

```

#categorize function
categorize = function (x, quantiles=(1:3)/4) {

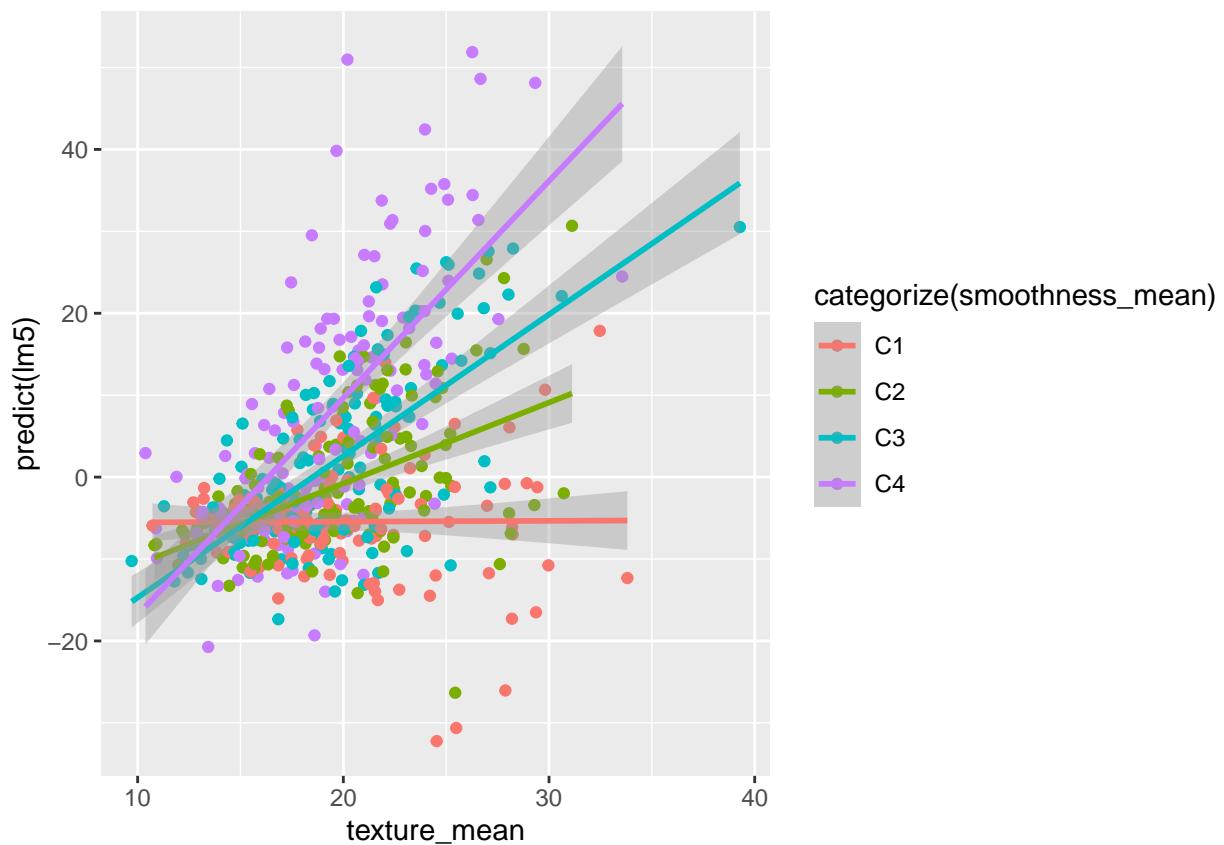
```

```

cutoffs = quantile (x, quantiles)
n.cutoffs = length (cutoffs)
result = rep ("C1", length (x))
for (j in 1:n.cutoffs) {
  result [x > cutoffs [j]] = paste ("C", j+1, sep="")
}
return (result)
}
#Interaction plot between mean texture and mean smoothness
library (ggplot2)
qplot (texture_mean, predict (lm5), data= wisconsin, color=categorize (smoothness_mean)) +
  geom_smooth (method="lm")

## `geom_smooth()` using formula 'y ~ x'

```



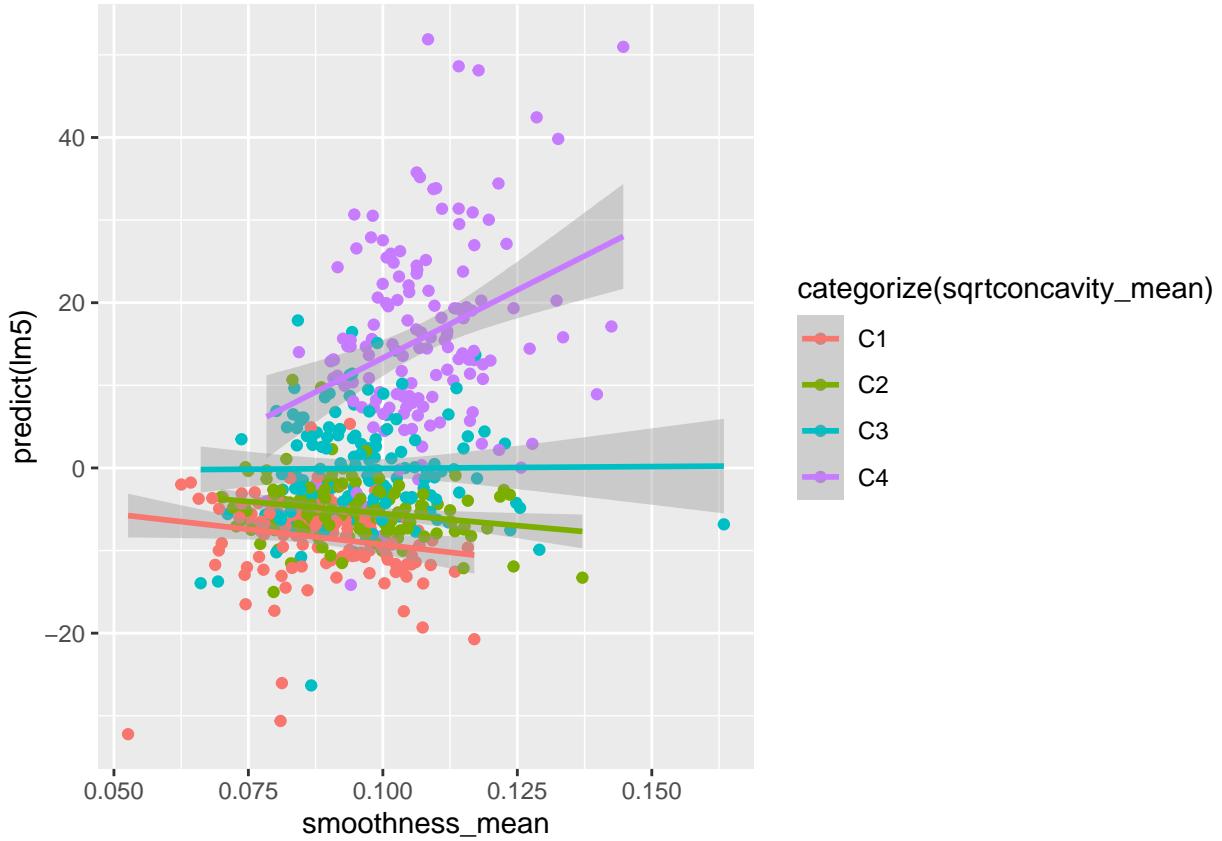
The interaction plot above shows that the relationship between the odds of a tumor being malignant and mean texture is stronger (steeper) for samples with higher mean smoothness values as opposed to lower smoothness values. This suggests that a sample's mean texture is more predictive of diagnosis (malignant vs benign) if on average it is more smooth.

```

# Interaction plot between mean smoothness and mean sqrt concavity
qplot (smoothness_mean, predict (lm5), data= wisconsin, color=categorize (sqrtconcavity_mean)) +
  geom_smooth (method="lm")

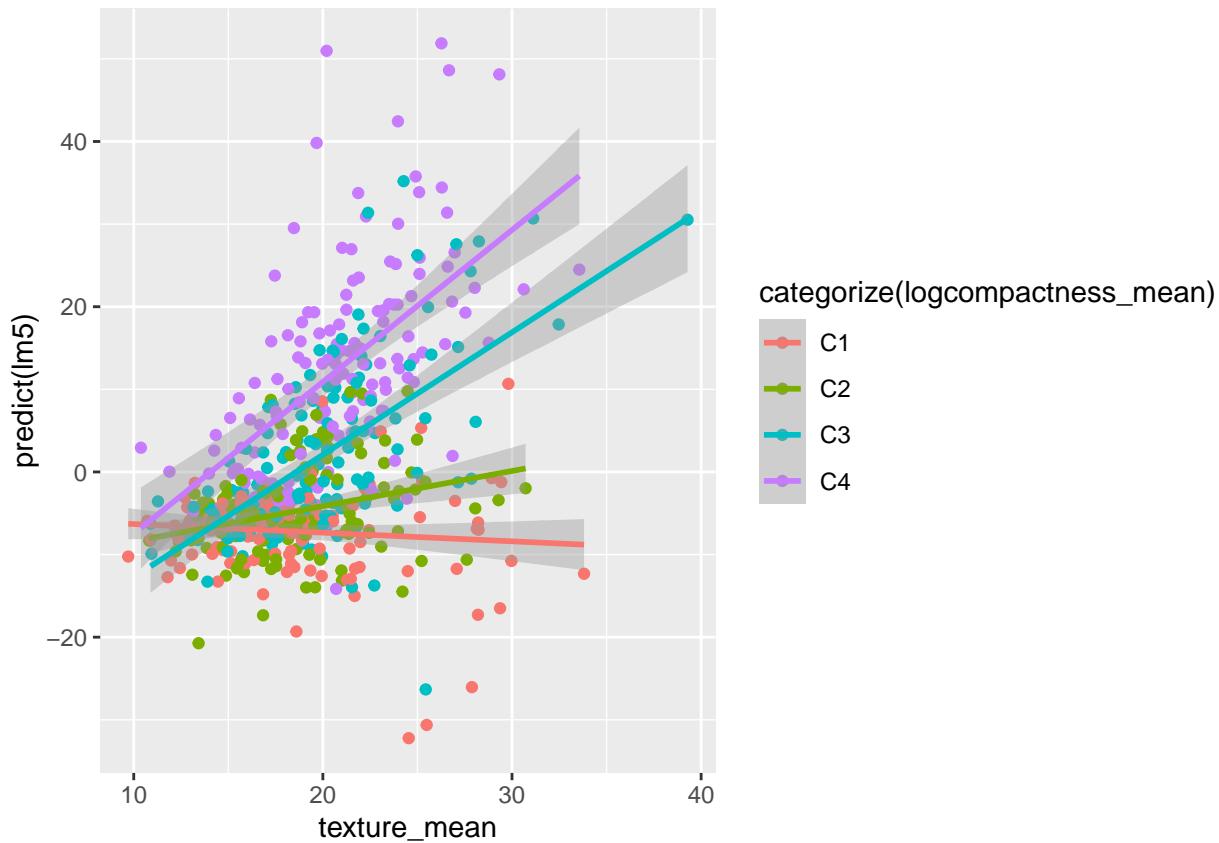
## `geom_smooth()` using formula 'y ~ x'

```



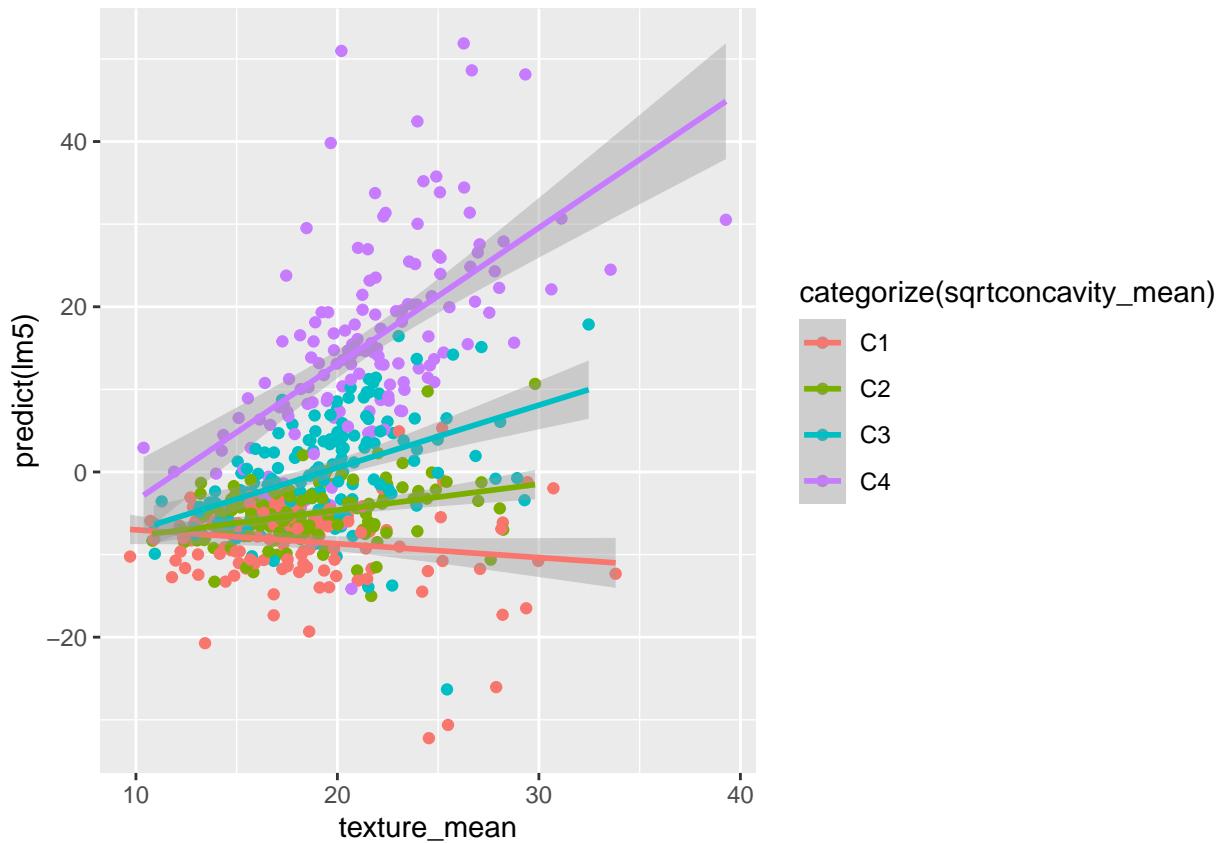
The interaction plot above shows that the relationship between mean smoothness and the odds of a tumor being malignant is minimal (close to zero) for cases in the lowest 75% of square root transformed mean concavity but strong for the cases with the highest square root transformed mean concavity. In other words, when mean square root transformed concavity is not high, mean smoothness is not very predictive of tumor diagnosis. However, when concavity is high, increasing mean smoothness is strongly associated with increased likelihood of a malignant tumor diagnosis.

```
# Interaction plot between mean texture and mean log compactness
qplot (texture_mean, predict (lm5), data= wisconsin, color=categorize (logcompactness_mean)) +
  geom_smooth (method="lm")
## `geom_smooth()` using formula 'y ~ x'
```



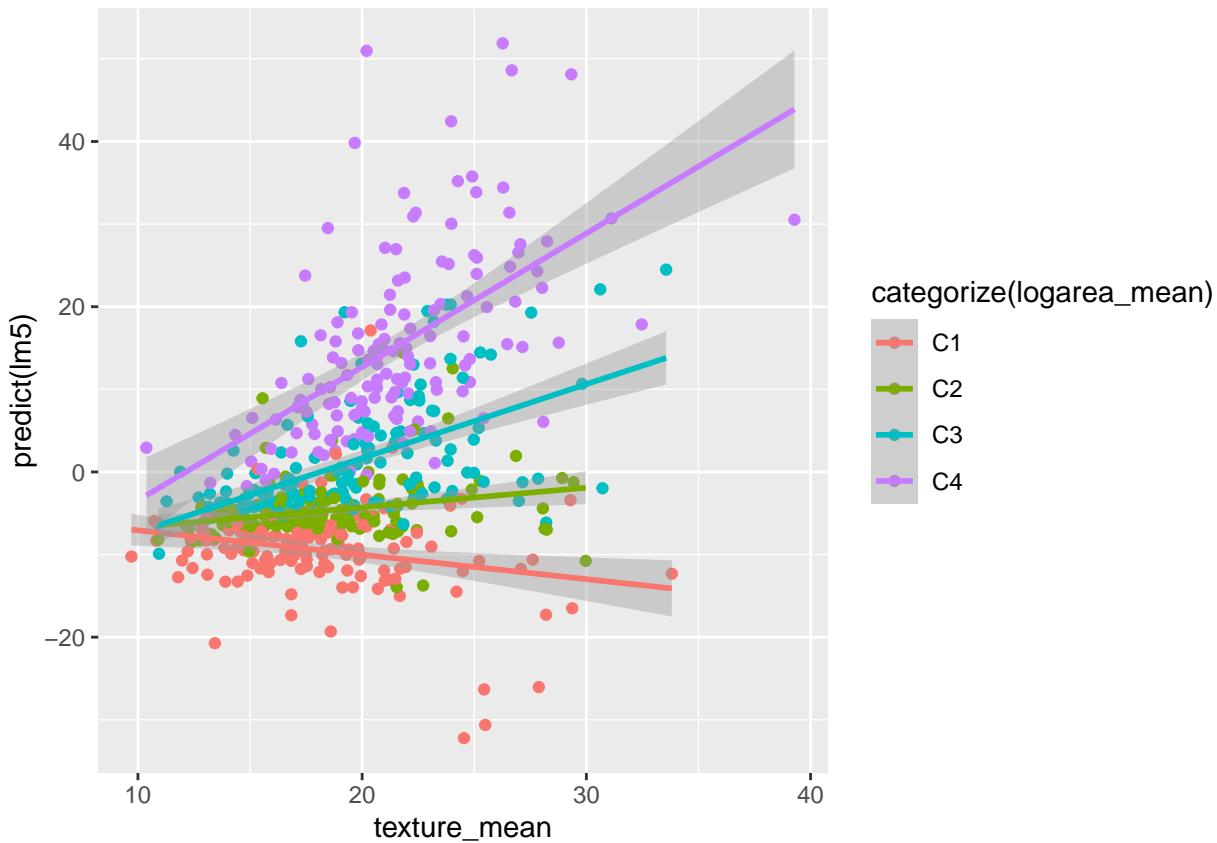
The interaction plot above shows that the relationship between the probability of a tumor being malignant and mean texture is stronger (steeper) for samples with higher mean log transformed compactness values as opposed to lower values. This suggests that a sample's mean texture is more predictive of tumor diagnosis if on average it is more compact.

```
# Interaction between mean texture and mean sqrt concavity
qplot (texture_mean, predict (lm5), data= wisconsin, color=categorize (sqrtconcavity_mean)) +
  geom_smooth (method="lm")
## `geom_smooth()` using formula 'y ~ x'
```



The interaction plot above shows that the relationship between the probability of a tumor being malignant and mean texture is stronger (steeper) for samples with higher mean square root transformed concavity values as opposed to lower values. This suggests that a sample's mean texture is more predictive of diagnosis (malignant vs benign) if on average it has greater concavity.

```
# Interaction between mean texture and mean log area
qplot (texture_mean, predict (lm5), data= wisconsin, color=categorize (logarea_mean)) +
  geom_smooth (method="lm")
## `geom_smooth()` using formula 'y ~ x'
```



The interaction plot above shows that for samples in the largest 50% of mean log area, the relationship between the probability of a tumor being malignant and mean texture is strong. Conversely, for samples with smaller mean log area values, the relationship between mean texture and the odds of a tumor being malignant is minimal. For cells with the smallest area, the relationship between odds of a malignant diagnosis and mean texture may become slightly negative which would suggest that in samples with low average cell area, increasing mean texture results in lower odds of a malignant diagnosis.

## Final Model

All predictors in the final model are included in the significant interaction effects making interpretations of the individual coefficients less meaningful. Our qualitative interpretations of the interaction effects can be found above.

Excluding interactions, the variables in our final model include mean texture, smoothness, square root concavity, log compactness, and log area.

A summary of the coefficient estimate, standard error, and p value for each variable:

- The mean texture coefficient estimate is 0.5738 (+/- 0.1047 SE) with a p value of 4.23e-08.
- The mean smoothness coefficient estimate is 206.496 (+/- 43.8422 SE) with a p value of 2.48e-06.
- The mean square root transformed concavity coefficient estimate is 29.9367 (+/- 6.2486 SE) with a p value of 1.66e-06.
- The mean log transformed compactness coefficient estimate is -4.1492 (+/- 1.5654 SE) with a p value of 0.008035.
- The mean log transformed area coefficient estimate is 12.2830 (+/- 1.8590 SE) with a p value of 3.91e-11.

The next table shows the odds ratios for the parameter estimates of the final model and their 95% confidence intervals:

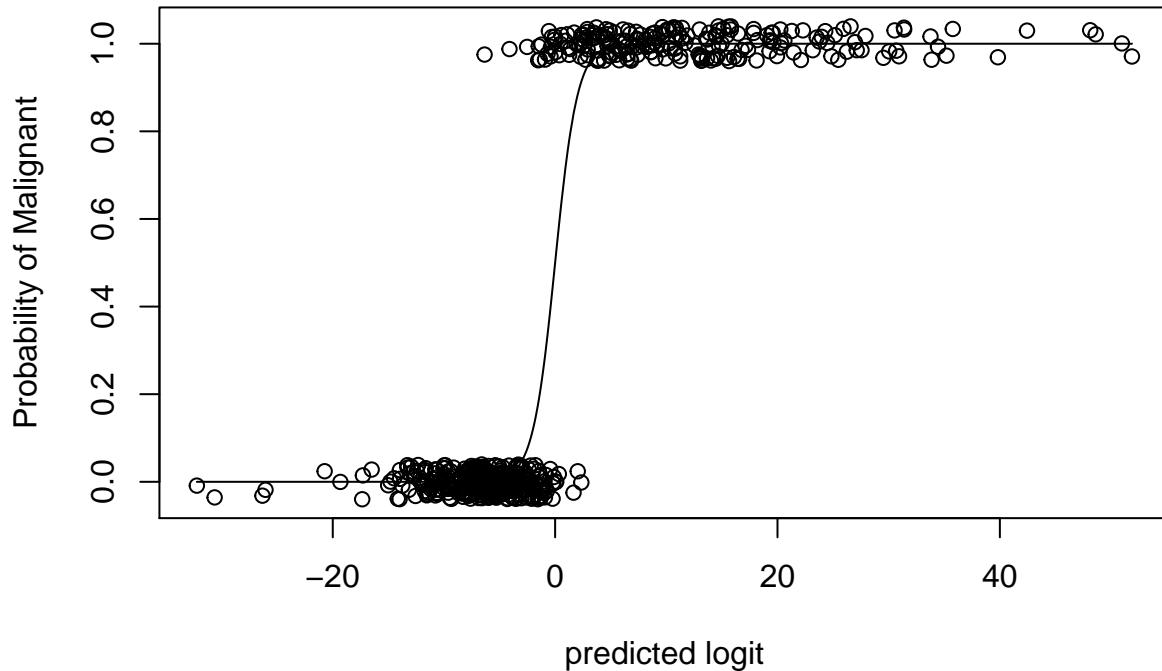












By observing the density of points when separated by our final model, there is a clear switch from benign to malignant tumor cells near the predicted logit of 0. The clear switch with little overlap from benign to malignant tumor cell responses as predicted logit increases suggests that our final model is well suited to accurately predict tumor diagnosis.

```
car::vif(lm5)
```

```
##          texture_mean.c           smoothness_mean.c
##                2.727401                4.862768
##          sqrtconcavity_mean.c      logcompactness_mean.c
##                5.365662                6.618617
##          logarea_mean.c   texture_mean.c:smoothness_mean.c
##                3.390512                4.073068
##  texture_mean.c:sqrtconcavity_mean.c  texture_mean.c:logcompactness_mean.c
##                5.396285                6.524692
##  texture_mean.c:logarea_mean.c smoothness_mean.c:sqrtconcavity_mean.c
##                1.795399                2.413216
```

After centering predictors and step wise regression, we have addressed most of our previous problems with collinearity. Mean log compactness and mean sqrt concavity and a subset of the interactions in which these predictors are included have VIF values just slightly larger than 5 (maximum VIF = 6.62). This suggests that we may have some moderate collinearity issues, but given that the VIF values are only marginally larger than the cutoff, we are confident proceeding with this model.

## Model Diagnostics

```
# Likelihood ratio test of final model
pchisq(lm5>null.deviance - lm5$deviance, lm5$df.null - lm5$df.residual, lower.tail = F)

## [1] 5.726579e-132

# Goodness of fit test of final model
1 - pchisq(lm5$deviance, lm5$df.residual)

## [1] 1
```

The likelihood ratio test shows that the as a whole has a significant effect on tumor diagnosis ( $p = 5.726579 \times 10^{-132}$ ). This suggests that the final model is useful for predicting diagnosis.

The goodness of fit test shows that the final model is not significantly different from a saturated model ( $p=1$ ). This suggests that our model is capturing most of the nuance in the association between the predictors and diagnosis and that there is no significant lack of fit.

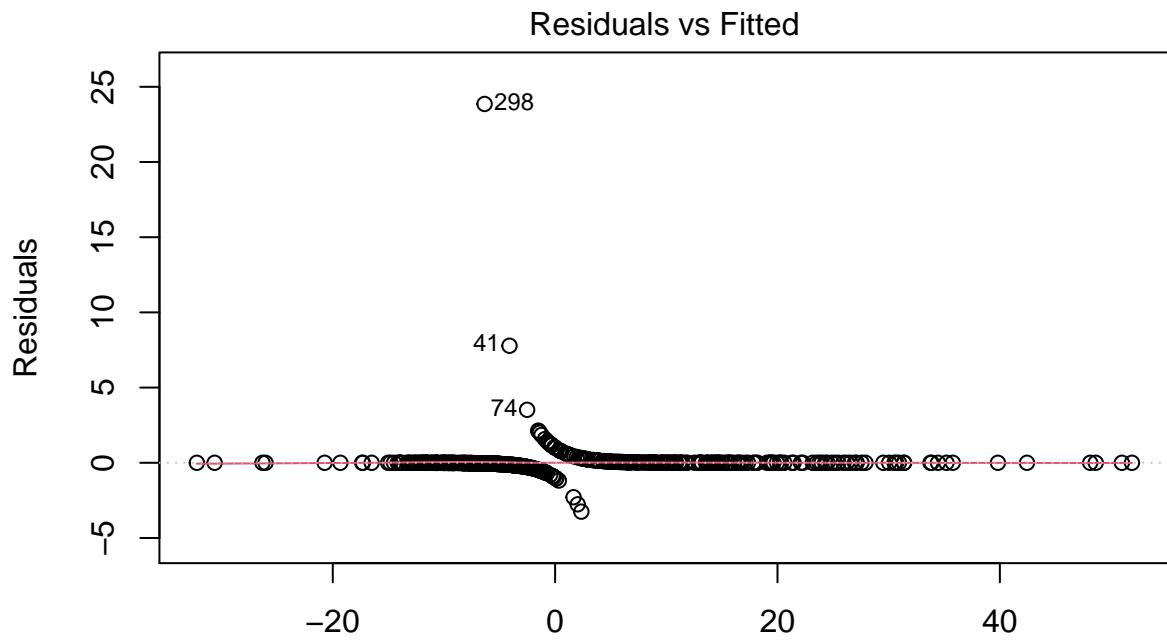
## Residual Plots and influence analysis

```
levcut = 2*(sqrt(11/569))
levcut

## [1] 0.2780803
```

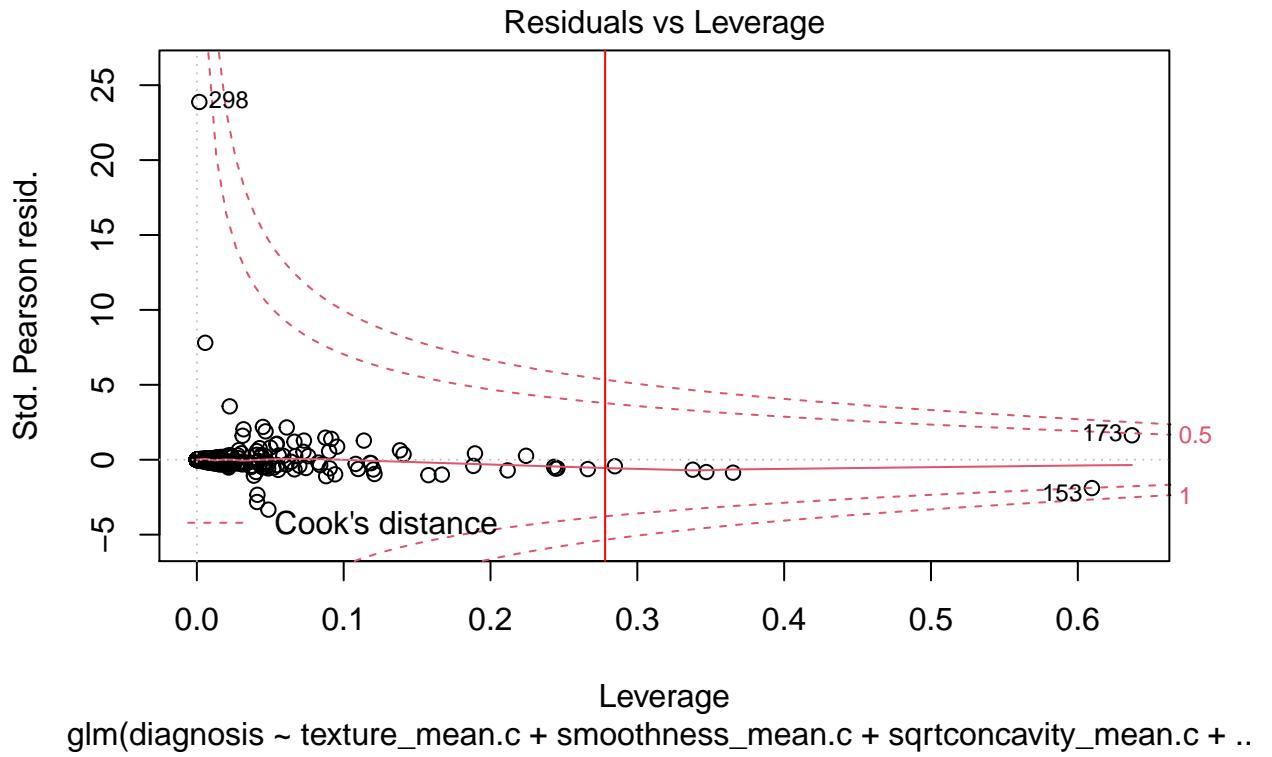
The leverage cut off is 0.278.

```
par(mfrow=c(1,1))
plot(lm5, which = c(1))
```



glm(diagnosis ~ texture\_mean.c + smoothness\_mean.c + sqrtconcavity\_mean.c + ..

```
plot(lm5, which = c(5))
abline(v=levcut, col="red")
```



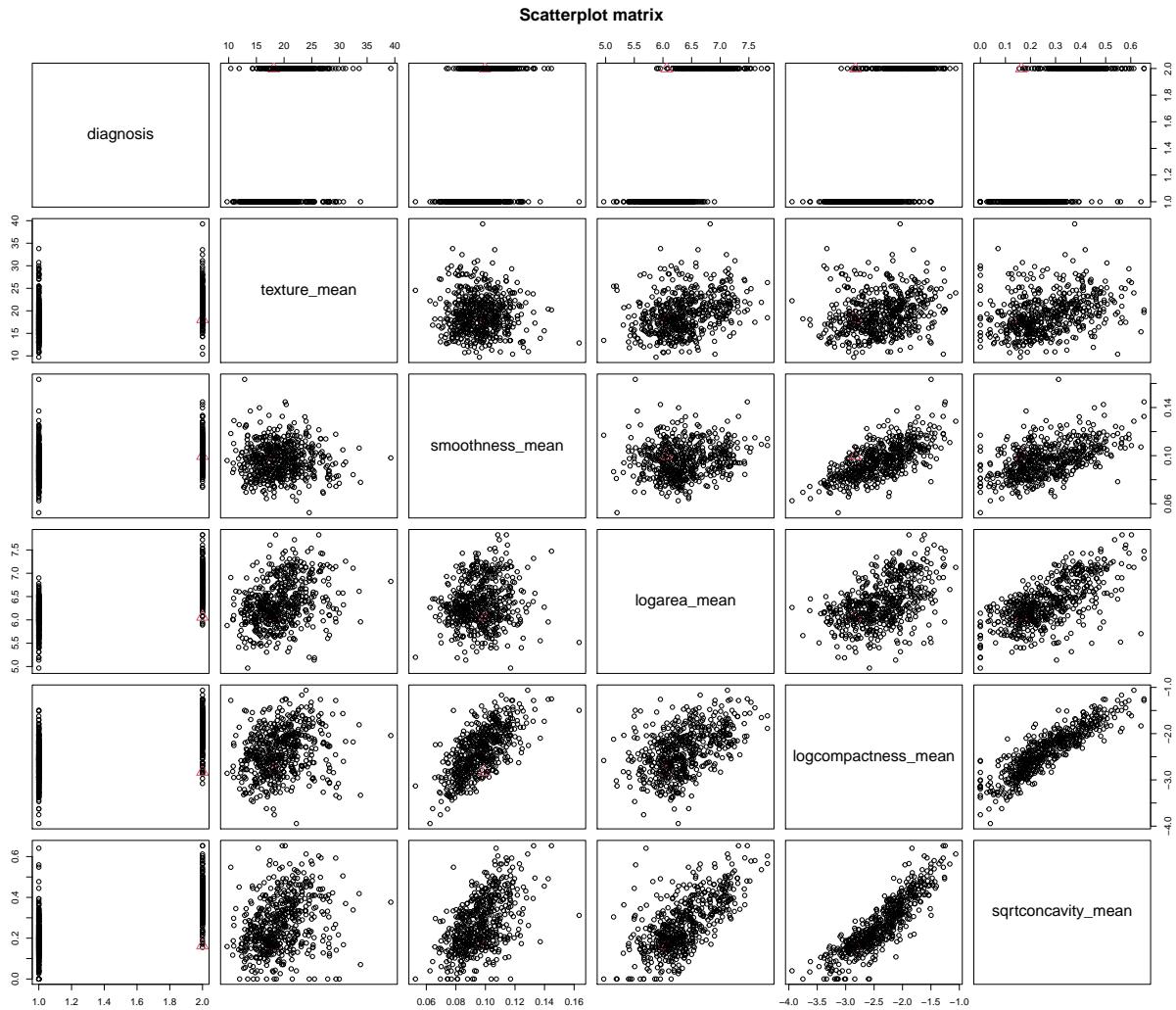
```
hatvals = hatvalues(lm5) > levcut
sum(hatvals)/569
```

```
## [1] 0.01054482
# summary of cooks distances
cooks <- as.data.frame(cooks.distance(lm5))
```

The smoothing spline in the residual plot for our first order model has little deviation from zero which suggests that our model has successfully met the assumption of linearity.

There are 6 samples (1.05% of the sample size) with leverage values above the cutoff. This is less than the expected 5% and suggests that the large majority of the samples in our sample have similar values for our predictor variables. Sample 153 and 173 have the highest leverage. While sample 173 is not quite influential, point 153 is a moderate influence point with a cooks distance value of .507. Points 41 and 298, while not influential due to low leverage, have very high positive residuals above the cut off of +/- 4 for standardized pearson residuals. Point 298 and 153 along with all influential points will be highlighted in a scatter plot matrix to examine the cause of their high leverage and residuals.

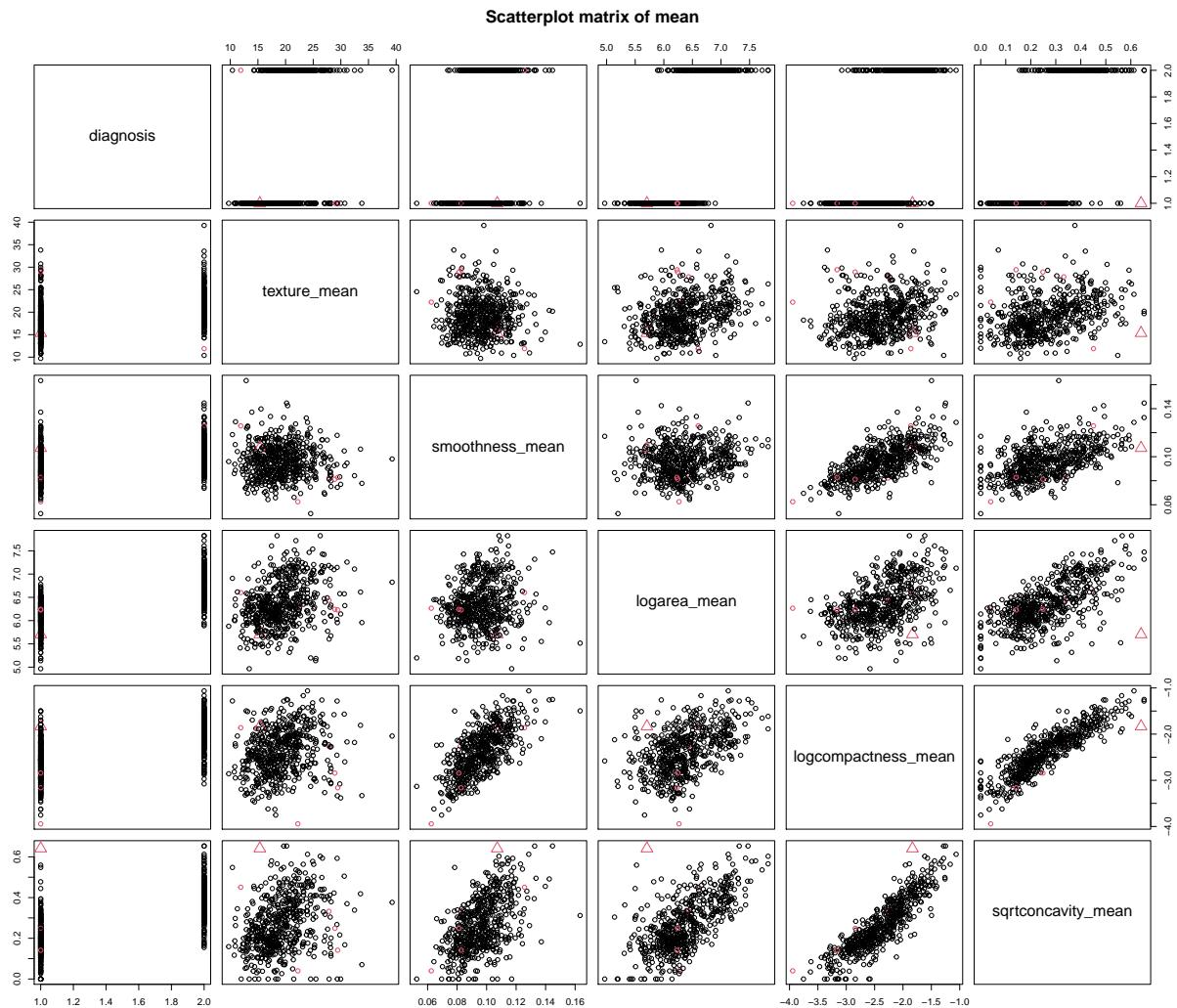
```
# Highlighting point 298 with a red triangle
pairs(wisconsin[, c(1, 3, 6, 12:14)], col= ifelse(1:569==298, 2, 1), cex=ifelse(1:569==298, 2, 1), pch=
```



As expected from the low influence value of sample 298, its values for each individual predictor variable are close to the mean value of all the cases. When separating cases by diagnosis however, sample 298's values for mean log area, log compactness, and square root concavity are towards the far left end of the distribution of malignant cases. These predictor variable values which are more associated with benign diagnosis rather than malignant diagnosis are likely the reason for the large standardized pearson residual. Because sample 298 is not influential and there is no clear evidence of data misentry or a plausible explanation for its odd pattern, we will leave sample 298 in our data set.

*# Highlighting point 153 as a triangle and coloring all high leverage points*

```
pairs(wisconsin[, c(1, 3, 6, 12:14)], col= (hatvals > levcut) + 1, cex=ifelse(1:569==153, 2, 1), pch=if
```



It appears that sample 153, one of our highest leverage points and the only moderately influential point in our data set, is explained primarily by its mean square root concavity value. Sample 153 has a mean square root concavity value that is towards the far right side of the distribution of this predictor for all cases making it the primary contributor to the samples high leverage. Additionally, when looking only at benign cases, sample 153's mean square root concavity measurement is even more abnormally large. It appears that this mean square root concavity value in addition to an abnormally large mean log compactness measurement for benign cases are the primary contributors to the samples residual magnitude. Because sample 153 is just barely moderately influential, we will keep it in our data set and proceed with the current model without any remedial measures.

## ROC Curve

```
par(mfrow= c(1, 1))
library(ROCR)

## Warning: package 'ROCR' was built under R version 4.0.5
pred1 <- prediction(lm5$fitted.values, lm5$y)
perf1 <- performance(pred1,"tpr","fpr")
```

```

auc1 <- performance(pred1, "auc")@y.values[[1]]
auc1

## [1] 0.9928122

plot(perf1, lwd=2, col=2)
abline(0,1)
legend(0.6, 0.3, c(paste("AUC=", round(auc1,4), sep="")), lwd = 2, col=2)
roc.x = slot(perf1, "x.values") [[1]]
roc.y = slot(perf1, "y.values") [[1]]
cutoffs = slot(perf1, "alpha.values") [[1]]

auc.table = cbind.data.frame(cutoff = pred1@cutoffs, tp=pred1@tp, fp=pred1@fp, tn=pred1@tn, fn=pred1@fn)

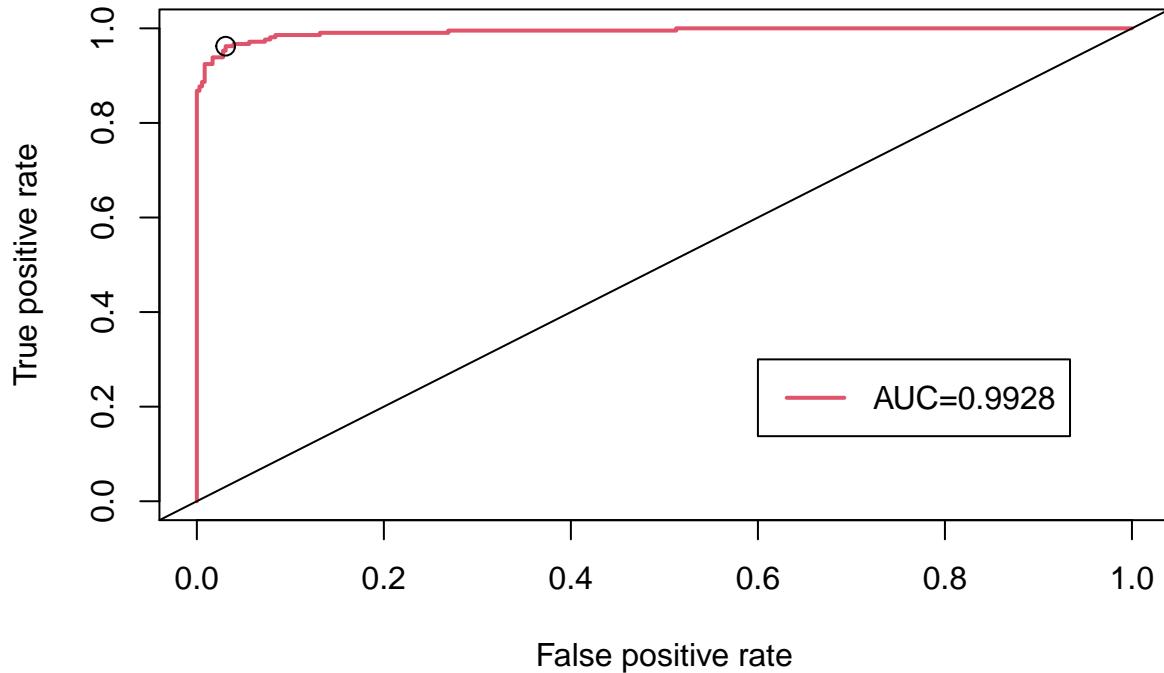
names(auc.table) = c("Cutoff", "TP", "FP", "TN", "FN")
auc.table$sensitivity = auc.table$TP / (auc.table$TP + auc.table$FN)
auc.table$specificity = auc.table$TN / (auc.table$TN + auc.table$FP)
auc.table$FalsePosRate = 1 - auc.table$specificity
auc.table$sens_spec = auc.table$sensitivity + auc.table$specificity
auc.table$ppv = auc.table$TP / (auc.table$TP + auc.table$FP)
auc.table$npv = auc.table$TN / (auc.table$TN + auc.table$FN)

auc.best = auc.table[auc.table$sens_spec == max(auc.table$sens_spec),]
auc.best

##      Cutoff TP FP TN FN sensitivity specificity FalsePosRate sens_spec
## 185 0.3361295 204 11 346 8    0.9622642   0.9691877   0.03081232 1.931452
##          ppv      npv
## 185 0.9488372 0.9774011

points(auc.best$FalsePosRate, auc.best$sensitivity, cex=1.3)

```



Our concordance index is 0.9928, meaning that the model is highly effective at predicting tumor diagnosis. Our optimal cutoff for determining whether a case is benign or malignant was 0.3361. This value maximized sensitivity and specificity. At this cut off the true positive rate or sensitivity was 0.962 and the true negative rate or specificity was 0.969. The positive predictive value or the probability that cases with a malignant prediction are truly malignant is 0.949. The negative predictive value or the probability that cases with a prediction are truly benign is 0.977. Again, these values suggest that our model is highly effective at predicting tumor diagnosis with minimal false predictions.

## Predictions using the final model

```

preds = predict(lm5, se.fit = T)
pred.df = cbind.data.frame(wisconsin, as.data.frame(preds))
pred.df$lwr = pred.df$fit - 1.96 * pred.df$se.fit
pred.df$upr = pred.df$fit + 1.96 * pred.df$se.fit
pred.df$fit.pr = round(exp(pred.df$fit) / (1 + exp(pred.df$fit)), 3)
pred.df$lwr.pr = round(exp(pred.df$lwr) / (1 + exp(pred.df$lwr)), 3)
pred.df$upr.pr = round(exp(pred.df$upr) / (1 + exp(pred.df$upr)), 3)

pred.df [c(2, 10, 298, 38, 68, 129), c(1, 16:17, 19, 21:22, 28:30)]
```

	diagnosis	texture_mean.c	smoothness_mean.c	sqrtconcavity_mean.c
## 2	M	-1.5196485	-0.011620281	0.02743367
## 10	M	4.7503515	0.022239719	0.20940551
## 298	M	-1.1496485	0.003319719	-0.10349469
## 38	B	-0.8696485	-0.006530281	-0.10729190

```

## 68      B     -0.2496485    -0.014970281    -0.07476675
## 129     B     -2.8996485     0.018639719     0.06998816
## logcompactness_mean.c logarea_mean.c fit.pr lwr.pr upr.pr
## 2          -0.1623572    0.8267372   0.997   0.968   1.000
## 10         0.9517332   -0.1979772   1.000   0.999   1.000
## 298        -0.4473302   -0.2968448   0.002   0.000   0.012
## 38          -0.8986392   -0.1020750   0.032   0.005   0.181
## 68          -0.6768774   -0.3865802   0.001   0.000   0.005
## 129        0.6696005    0.1507867   0.913   0.694   0.980

```

In order to test how effective our final model is, we made 3 predictions of malignant samples and 3 predictions of benign samples. Our prediction cut off as determined by AUC, which maximizes sensitivity and specificity is 0.3361. Any cases with probabilities above this threshold are predicted to be malignant whereas cases with probabilities below this threshold are predicted to be benign. To challenge our model we included in our predictions sample 298, our highest positive residual, and sample 129, our largest negative residual. Our model correctly predicted 4 of these 6 cases.

Outside of case 298 and 129, our model correctly predicted the other cases relatively easily given that the prediction interval for these predictions did not come close to overlapping with the established cut off.

```

pred.df$pred.fail = ifelse (pred.df$fit.pr >= auc.best$Cutoff[1], "Pred.Yes",
"Pred.No")
table (wisconsin$diagnosis, pred.df$pred.fail)

```

```

##
##      Pred.No Pred.Yes
##      B      346      11
##      M       9      203

```

In order to situate our 6 predictions within the context of the full data set we revisit the contingency table. Altogether, our model correctly predicted 96.2% of malignant cases and 96.9% of benign cases suggesting that our 2 incorrect predictions were representative of a very small proportion of incorrectly identified cases in the entire data set. Overall, our model is very accurate and highly effective for predicting tumor diagnosis with very few false predictions.

## Conclusions

Altogether, our model shows that mean texture, smoothness, square root concavity, log compactness, and log area are significant predictors of tumor diagnosis. Between these variables there were significant interactions between texture and smoothness, square root concavity, log compactness, and log area as well as an interaction between smoothness and square root concavity. The AUC of this model was .9928 which shows that this combination of predictors is very effective at predicting tumor diagnosis. A high p value ( $p=1$ ) for a goodness of fit test of the model shows that there is no significant lack of fit. The accuracy of the model at predicting both benign and malignant tumors is shown by the sensitivity and specificity at the established cut off. Altogether, 96.2% of malignant cases and 96.9% of benign cases were correctly predicted by our model.

What questions are not answered by our model?

-We wonder how well our our model would perform given new data. It would be interesting to conduct cross validation to test this.

-Additionally, we are curious to learn more about how models like this one are applied in health care. Is it good enough to stand-in for an actual doctor making a diagnosis? Could these models be applicable to other kinds of cancers as well?

-It would be interesting to better understand the biological basis for the trends we observed in predictor variables, especially in the interaction effects. For instance, why is texture strongly associated with breast cancer when cells are relatively large, but not when cells are smaller?

- An explanation for the abnormal pattern of sample 298 could lead to interesting/important scientific findings.