

Initial Multiple Regression Analysis of Color Values in Relation to Mean Area for Cancer Cells

Caleb Munson, Luke Geel, Leah Dion, Everett Henderson, Jiun Tseng

April 2021

Description of initial multiple linear regression model

Our initial model includes the 7 predictors (respectively Mean Area of Red, Green, Blue, HSV, LAB, HE, BR) in relation to the response, mean area across each nuclei. Each predictor is accounted for linearly. b_0 is the expected mean area assuming each predictor is zero.

$$\hat{Y} = 183.3774 - 5.8856X_1 - 4.3324X_2 - 2.7434X_3 + 12.7141X_4 + 4.1442X_5 + 1.1101X_6 - 0.1687X_7$$

Summary of estimated regression coefficients and standard errors

Summary of b_i values:

- $b_0 = 183.3774$
- $b_1 = -5.8856$
- $b_2 = -4.3324$
- $b_3 = -2.7434$
- $b_4 = 12.7141$
- $b_5 = 4.1442$
- $b_6 = 1.1101$
- $b_7 = -0.1687$

Summary of the absolute value of the residuals:

- mean = 90.691
- s.d. = 81.0104
- min = 2.099
- max = 585.026
- median = 72.985

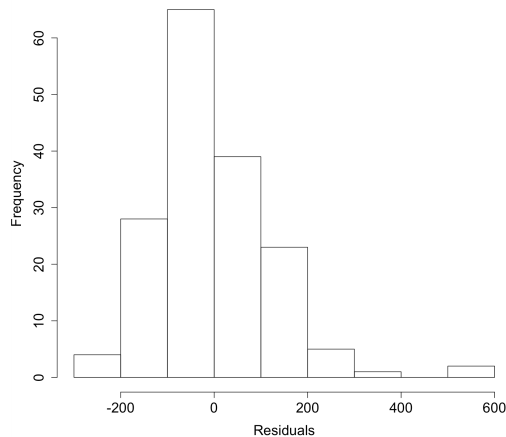


Figure 1: Histogram of the residuals approximately resembles a normal distribution.

Graphical assessment of residuals and their variance

Figure 2: Standard errors against fitted values.

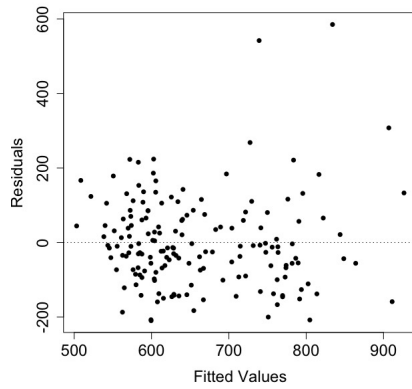
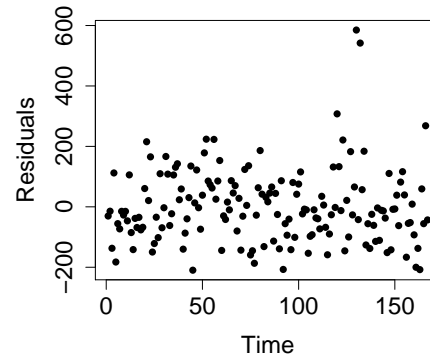


Figure 3: Time plot of the residuals.



Plots of residuals against each predictor

Figure 4: Residuals against mean red intensity.

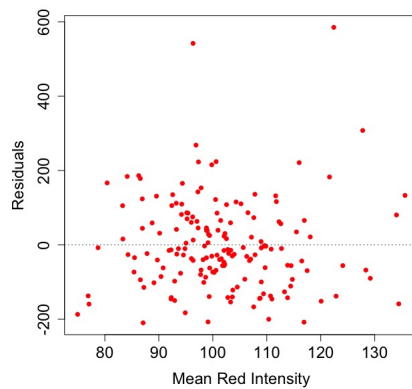


Figure 5: Residuals against mean blue intensity

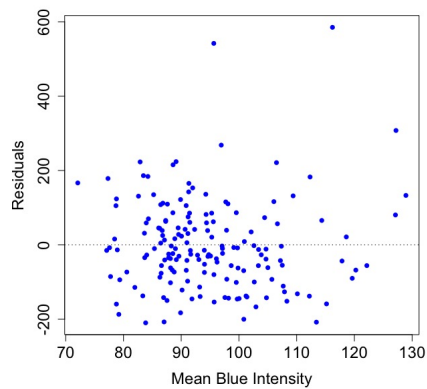


Figure 6: Residuals against mean green intensity.

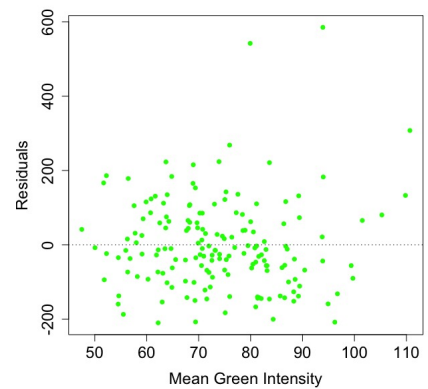


Figure 7: Residuals against mean HSV intensity.

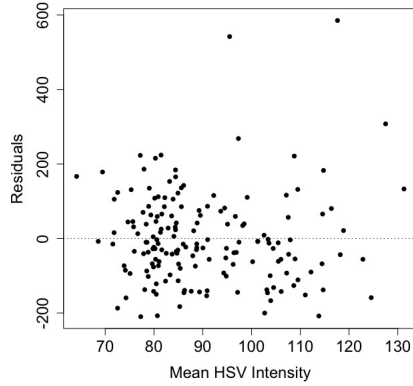


Figure 8: Residuals against mean Lab intensity.

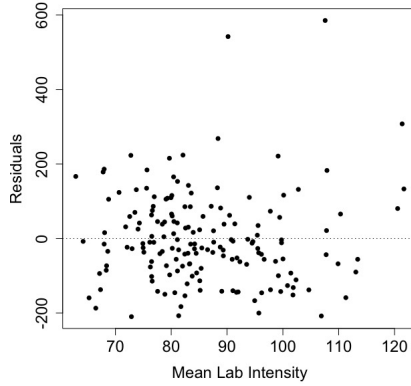


Figure 9: Residuals against mean HE intensity.

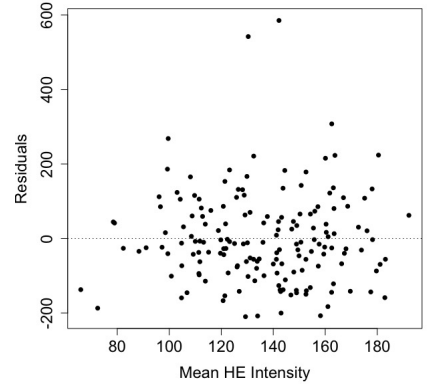
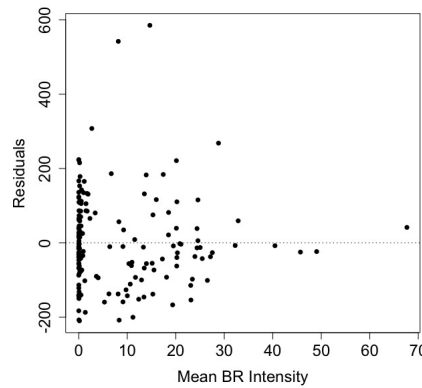


Figure 10: Residuals against mean BR intensity.



Evidence of violations of linear model assumptions

There are 2 extreme values correlating to nuclei with ID 130 and 132. A possible transformation would be to exclude these outlier values. These values severely skewed the residuals in the positive direction. In addition, since the summation of errors equal 0, we predict that the regression function is linear and the error terms appear to be independent. In the BR residual plot, we should note that many of the BR values are 0, and that the residuals are spread less the higher the BR intensity is. This may mean that we want to transform the BR variable in some way.

Key takeaways

Examining the residual plots, it looks like the multiple variable regression fits the model fairly well overall. The graphs of mean red, green, and blue intensity vs residuals are very similar to each other, suggesting that each one does equally well at predicting the response mean area. Additionally, we believe that each predictor works equally well at any frequency, so how well the model works doesn't change with frequency. The mean BR data was noticeably different than the other variables as the mean BR data was skewed more towards 0 than any other variable, which we will pay attention to upon further statistical analysis.

```

2 cancer <- read.csv("/Users/leahleahleah/Downloads/cancer.csv")
3
4 Y <- cancer$Mean_Area
5 X1 <- cancer$Mean_mean_R
6 X2 <- cancer$Mean_mean_G
7 X3 <- cancer$Mean_mean_B
8 X4 <- cancer$Mean_mean_HSV
9 X5 <- cancer$Mean_mean_Lab
10 X6 <- cancer$Mean_mean_HE
11 X7 <- cancer$Mean_mean_BR
12
13 lm.canc <- lm(Y ~ X1 + X2 + X3 + X4 + X5 + X6 + X7)
14
15 #extract b-values
16 b0 <- summary(lm.canc)$coef[1,1]
17 b1 <- summary(lm.canc)$coef[2,1]
18 b2 <- summary(lm.canc)$coef[3,1]
19 b3 <- summary(lm.canc)$coef[4,1]
20 b4 <- summary(lm.canc)$coef[5,1]
21 b5 <- summary(lm.canc)$coef[6,1]
22 b6 <- summary(lm.canc)$coef[7,1]
23 b7 <- summary(lm.canc)$coef[8,1]
24
25 #extract residuals
26 e <- lm.canc$residuals
27 sum(abs(e))
28 sd(abs(e))
29 #extract fitted values
30 Y.hat <- lm.canc$fitted.values
31 #histogram of residuals
32 hist(e, xlab = "Residuals", ylab = "Frequency", cex.lab = 1.5,
33      cex.axis = 1.5)
34
35 #plot of residuals vs fitted values
36 plot(Y.hat, e, pch = 16,
37      xlab = "Fitted Values", ylab = "Residuals", cex.lab = 1.5,
38      cex.axis = 1.5)
39 abline(h = 0, lty = 3)
40
41 #plot of residuals vs mean red intensity
42 plot(X1, e, pch = 16, col = 'red',
43      xlab = "Mean Red Intensity", ylab = "Residuals", cex.lab =
44      cex.axis = 1.5)
45 abline(h = 0, lty = 3)
46
47 #plot of residuals vs mean green intensity
48 plot(X2, e, pch = 16, col = 'green',
49      xlab = "Mean Green Intensity", ylab = "Residuals", cex.lab
50      cex.axis = 1.5)
51 abline(h = 0, lty = 3)

```

```

53 #plot of residuals vs mean blue intensity
54 plot(X3, e, pch = 16, col = 'blue',
55      xlab = "Mean Blue Intensity", ylab = "Residuals", cex.lab = 1.5,
56      cex.axis = 1.5)
57 abline(h = 0, lty = 3)
58
59 #plot of residuals vs mean hsv intensity
60 plot(X4, e, pch = 16,
61      xlab = "Mean HSV Intensity", ylab = "Residuals", cex.lab = 1.5,
62      cex.axis = 1.5)
63 abline(h = 0, lty = 3)
64
65 ##plot of residuals vs mean lab intensity
66 plot(X5, e, pch = 16,
67      xlab = "Mean Lab Intensity", ylab = "Residuals", cex.lab = 1.5,
68      cex.axis = 1.5)
69 abline(h = 0, lty = 3)
70
71 #plot of residuals vs mean he intensity
72 plot(X6, e, pch = 16,
73      xlab = "Mean HE Intensity", ylab = "Residuals", cex.lab = 1.5,
74      cex.axis = 1.5)
75 abline(h = 0, lty = 3)
76
77 #plot of residuals vs mean br intensity
78 plot(X7, e, pch = 16,
79      xlab = "Mean BR Intensity", ylab = "Residuals", cex.lab = 1.5,
80      cex.axis = 1.5)
81 abline(h = 0, lty = 3)

```