

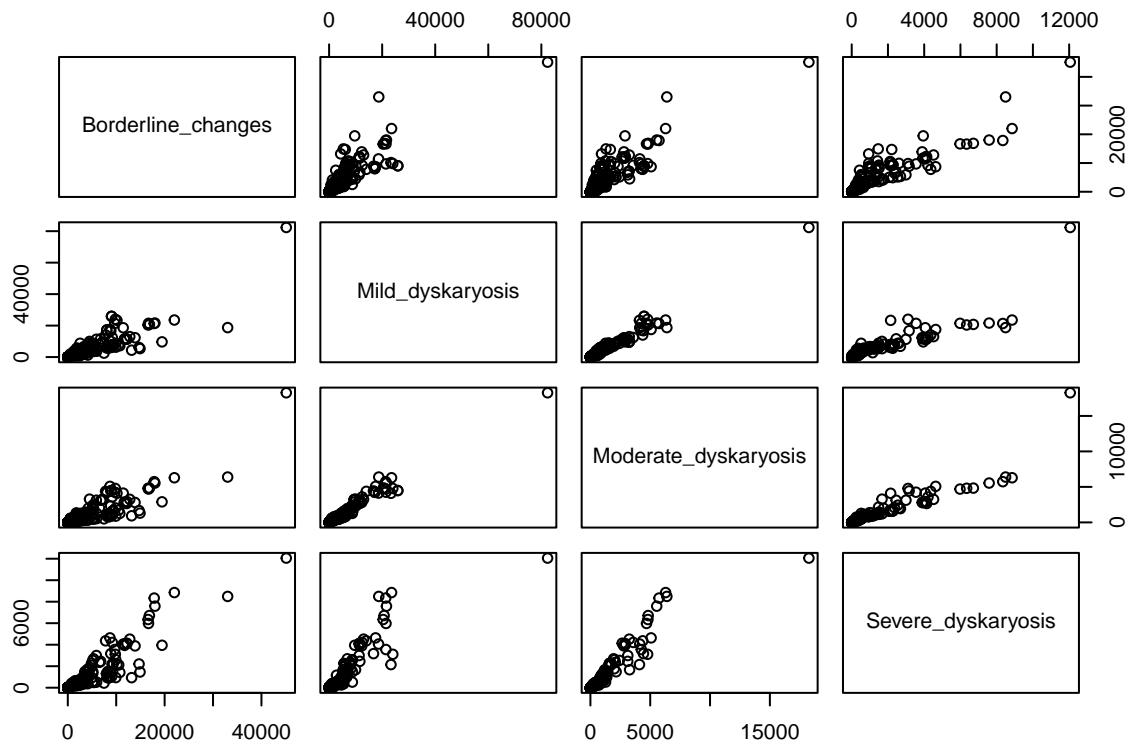
The above was all copied from the project check in 1 document. It will not appear in the knitted file because it is suppressed.

First, we will subset the data for borderline changes, mild, moderate, and severe dyskaryosis into a separate data frame. Then the pairs function will be run to determine any reasonable relationships.

```
## Subset source data for only borderline changes and mild, moderate, severe dyskaryosis
```

```
comparison_data <- source[,c("Borderline_changes", "Mild_dyskaryosis", "Moderate_dyskaryosis", "Severe_dyskaryosis")]
comparison_data$Severe_dyskaryosis <- as.numeric(comparison_data$Severe_dyskaryosis)
```

```
pairs(comparison_data)
```



First, we will run a linear regression for all of these relationships as they appear to have decently strong relationships.

### Unscaled plots

```
# Mild Dyskaryosis
```

```
lm_model_mild <- lm(Mild_dyskaryosis ~ Borderline_changes, data = comparison_data)
summary(lm_model_mild)
```

```
##
```

```
## Call:
```

```
## lm(formula = Mild_dyskaryosis ~ Borderline_changes, data = comparison_data)
```

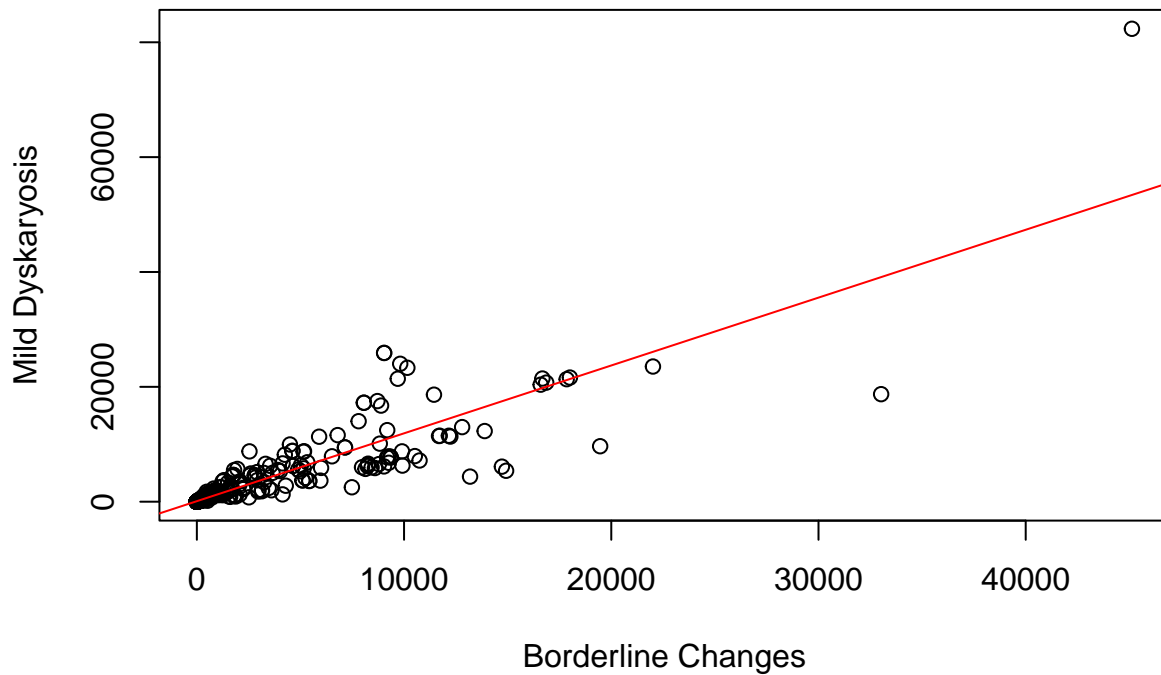
```
##
```

```
## Residuals:
```

```
##      Min      1Q   Median      3Q      Max
## -20385.1  -115.6   -84.5   190.5 28938.7
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    88.63154   210.22556    0.422   0.674
## Borderline_changes  1.18161    0.03679   32.117  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3338 on 328 degrees of freedom
## Multiple R-squared:  0.7587, Adjusted R-squared:  0.758
## F-statistic: 1031 on 1 and 328 DF, p-value: < 2.2e-16
```

```
plot(comparison_data$Borderline_changes, comparison_data$Mild_dyskaryosis,
     main = "Mild Dyskaryosis vs Borderline Changes",
     xlab = "Borderline Changes",
     ylab = "Mild Dyskaryosis")
abline(lm_model_mild, col = 'red')
```

## Mild Dyskaryosis vs Borderline Changes

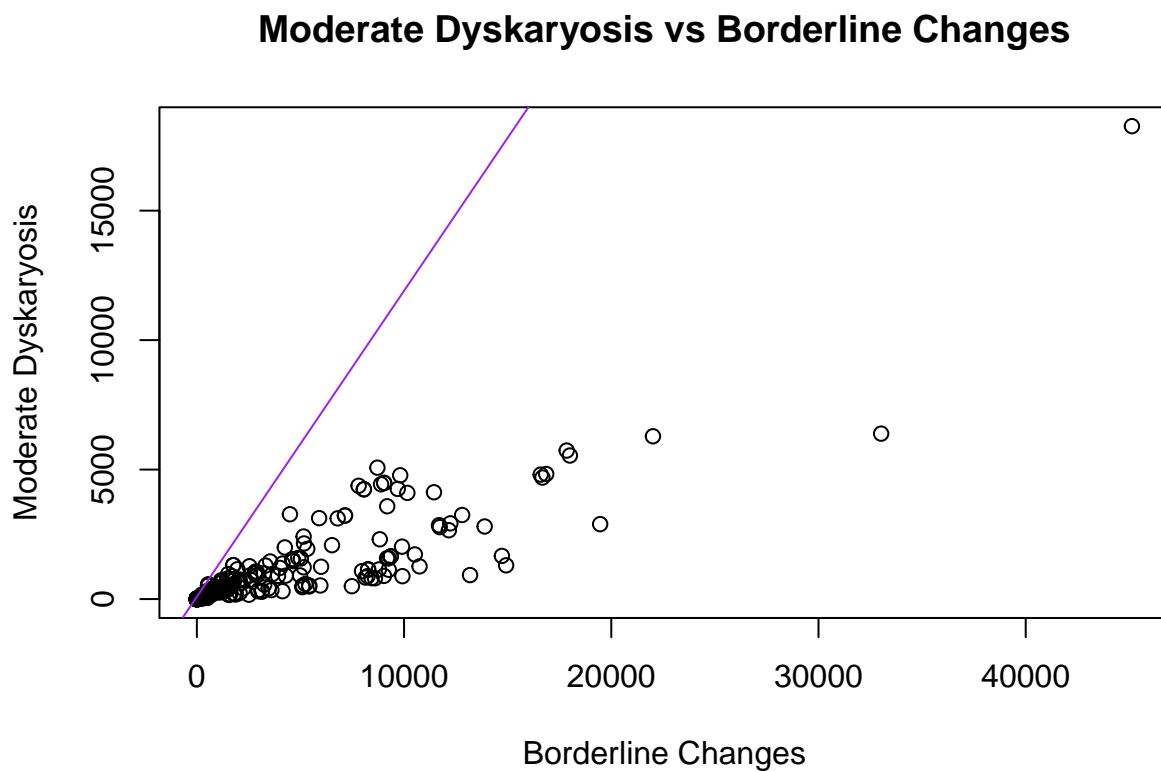


```
# Moderate Dyskaryosis
lm_model_mod <- lm(Moderate_dyskaryosis ~ Borderline_changes, data = comparison_data)
summary(lm_model_mod)
```

```
##
```

```
## Call:
## lm(formula = Moderate_dyskaryosis ~ Borderline_changes, data = comparison_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2835.2    -1.0      6.0     80.4   5741.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -5.414850  46.102620  -0.117   0.907
## Borderline_changes  0.277578   0.008068  34.403 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 732.1 on 328 degrees of freedom
## Multiple R-squared:  0.783, Adjusted R-squared:  0.7823
## F-statistic: 1184 on 1 and 328 DF, p-value: < 2.2e-16
```

```
plot(comparison_data$Borderline_changes, comparison_data$Moderate_dyskaryosis,
     main = "Moderate Dyskaryosis vs Borderline Changes",
     xlab = "Borderline Changes",
     ylab = "Moderate Dyskaryosis")
abline(lm_model_mild, col = 'purple')
```



```
# Severe Dyskaryosis
lm_model_severe <- lm(Severe_dyskaryosis ~ Borderline_changes, data = comparison_data)
summary(lm_model_severe)
```

```
##
## Call:
## lm(formula = Severe_dyskaryosis ~ Borderline_changes, data = comparison_data)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
##	-2752.9	-48.3	51.9	61.2	3327.0

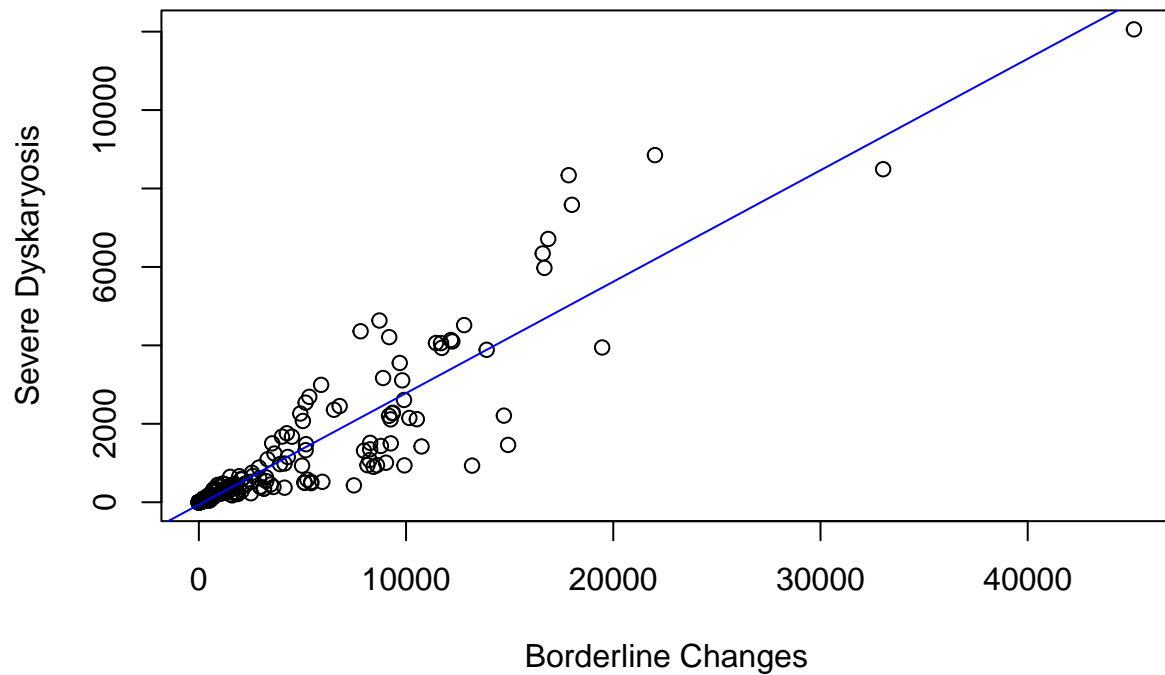
```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	-60.208265	52.390598	-1.149	0.252
## Borderline_changes	0.284191	0.008057	35.271	<2e-16 ***

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 692.2 on 238 degrees of freedom
## (90 observations deleted due to missingness)
## Multiple R-squared:  0.8394, Adjusted R-squared:  0.8387
## F-statistic: 1244 on 1 and 238 DF, p-value: < 2.2e-16
```

```
plot(comparison_data$Borderline_changes, comparison_data$Severe_dyskaryosis,
     main = "Severe Dyskaryosis vs Borderline Changes",
     xlab = "Borderline Changes",
     ylab = "Severe Dyskaryosis")
abline(lm_model_severe, col = 'blue')
```

## Severe Dyskaryosis vs Borderline Changes

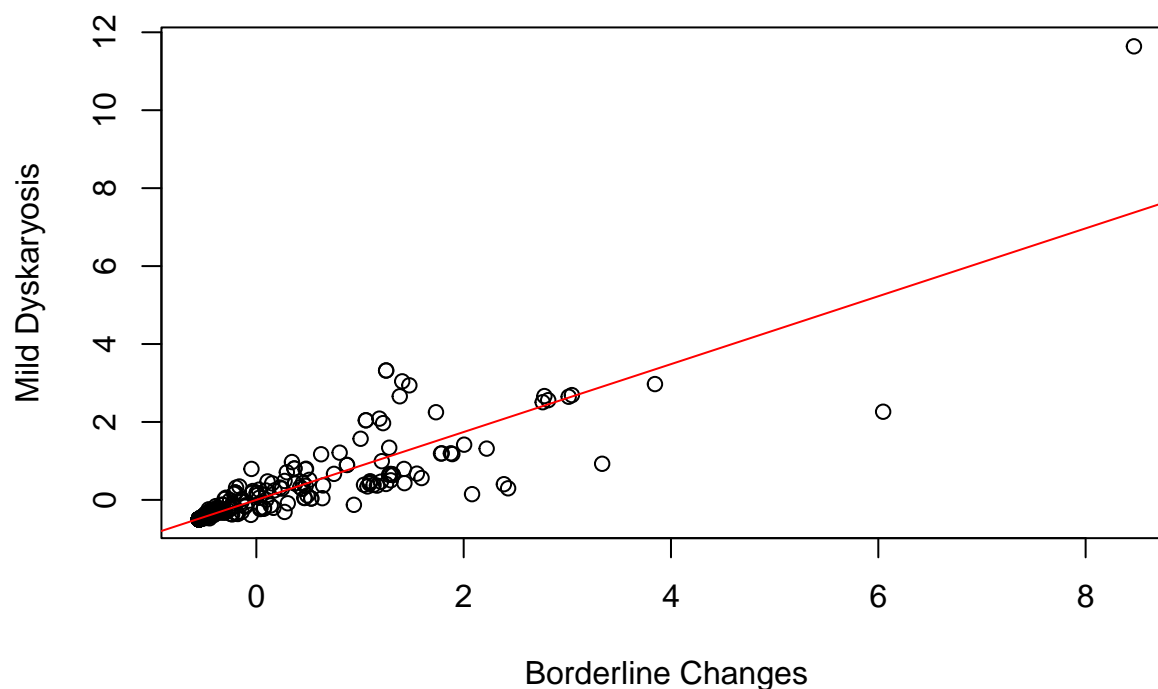


### Scaled Plots

```
scaled_data <- as.data.frame(scale(comparison_data))

# mild dyskaryosis scaled
scaled_mild_model <- lm(Mild_dyskaryosis ~ Borderline_changes, data = scaled_data)
plot(scaled_data$Borderline_changes, scaled_data$Mild_dyskaryosis,
     main = "Scaled Mild Dyskaryosis vs Borderline Changes",
     xlab = "Borderline Changes",
     ylab = "Mild Dyskaryosis")
abline(scaled_mild_model, col = 'red')
```

## Scaled Mild Dyskaryosis vs Borderline Changes



```
summary(scaled_mild_model)
```

```
##
## Call:
## lm(formula = Mild_dyskaryosis ~ Borderline_changes, data = scaled_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.0040 -0.0170 -0.0125  0.0281  4.2644
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.213e-16  2.708e-02   0.00      1
## Borderline_changes 8.711e-01  2.712e-02  32.12 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4919 on 328 degrees of freedom
## Multiple R-squared:  0.7587, Adjusted R-squared:  0.758
## F-statistic: 1031 on 1 and 328 DF, p-value: < 2.2e-16
```

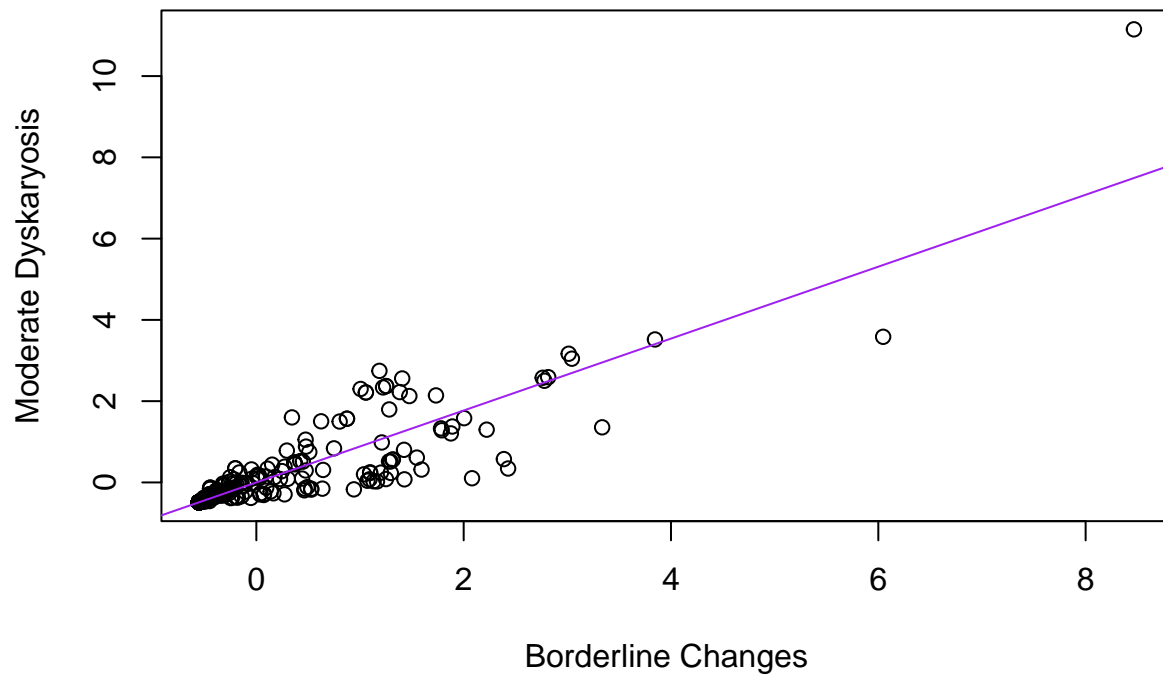
```
# moderate dyskaryosis scaled
scaled_mod_model <- lm(Moderate_dyskaryosis ~ Borderline_changes, data = scaled_data)
plot(scaled_data$Borderline_changes, scaled_data$Moderate_dyskaryosis,
     main = "Scaled Moderate Dyskaryosis vs Borderline Changes",
```

```

xlab = "Borderline Changes",
ylab = "Moderate Dyskaryosis")
abline(scaled_mod_model, col = 'purple')

```

## Scaled Moderate Dyskaryosis vs Borderline Changes



```
summary(scaled_mod_model)
```

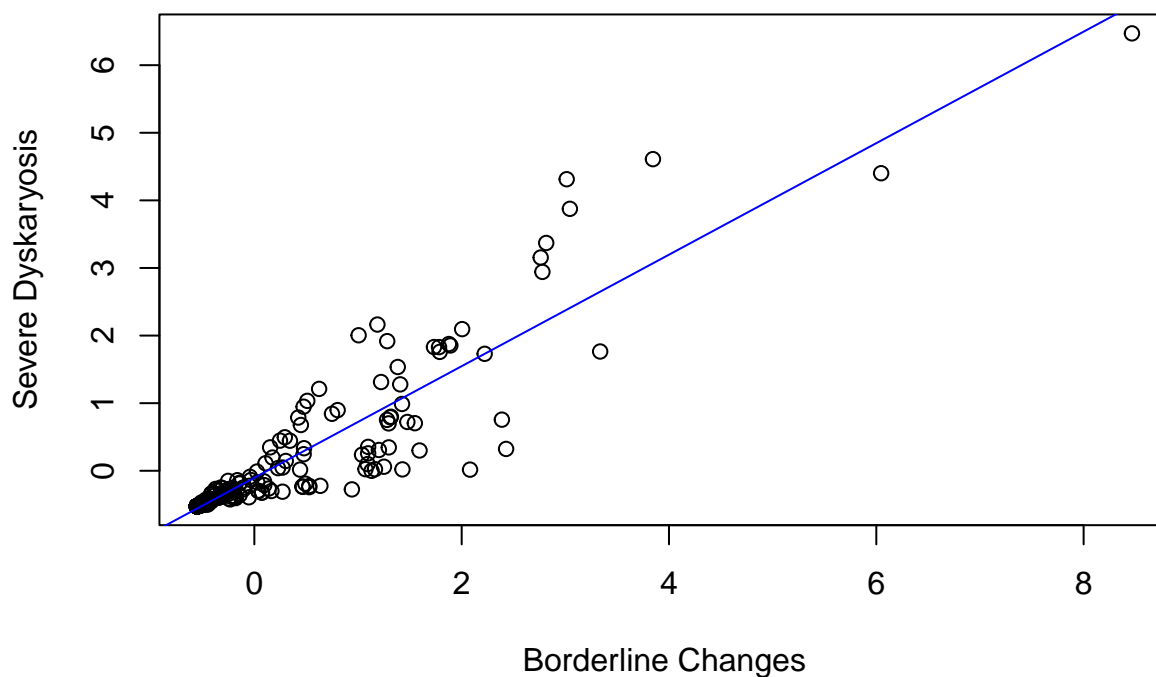
```

##
## Call:
## lm(formula = Moderate_dyskaryosis ~ Borderline_changes, data = scaled_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8067 -0.0007  0.0038  0.0512  3.6586
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.621e-16  2.568e-02   0.0      1
## Borderline_changes 8.849e-01  2.572e-02  34.4   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4665 on 328 degrees of freedom
## Multiple R-squared:  0.783, Adjusted R-squared:  0.7823
## F-statistic: 1184 on 1 and 328 DF, p-value: < 2.2e-16

```

```
# severe dyskaryosis scaled
scaled_severe_model <- lm(Severe_dyskaryosis ~ Borderline_changes, data = scaled_data)
plot(scaled_data$Borderline_changes, scaled_data$Severe_dyskaryosis,
     main = "Scaled Severe Dyskaryosis vs Borderline Changes",
     xlab = "Borderline Changes",
     ylab = "Severe Dyskaryosis")
abline(scaled_severe_model, col = 'blue')
```

## Scaled Severe Dyskaryosis vs Borderline Changes



```
summary(scaled_severe_model)
```

```
##
## Call:
## lm(formula = Severe_dyskaryosis ~ Borderline_changes, data = scaled_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.59706 -0.02801  0.03012  0.03551  1.93013
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -0.10225    0.02608   -3.92 0.000116 ***
## Borderline_changes  0.82478    0.02338   35.27 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```



```
## Residual standard error: 0.4016 on 238 degrees of freedom
## (90 observations deleted due to missingness)
## Multiple R-squared: 0.8394, Adjusted R-squared: 0.8387
## F-statistic: 1244 on 1 and 238 DF, p-value: < 2.2e-16
```

For the first linear regression model of Mild Dyskaryosis vs Borderline Changes, the coefficient of the slope of the line was significant with an alpha level of  $<0.001$ . The intercept of the line was not significant at any level. The adjusted R-squared value is 0.758, indicating that 75.8% of the variance in the data can be explained by the model. The F statistic is large with a p-value  $< 0.05$ , indicating that this regression is significant.

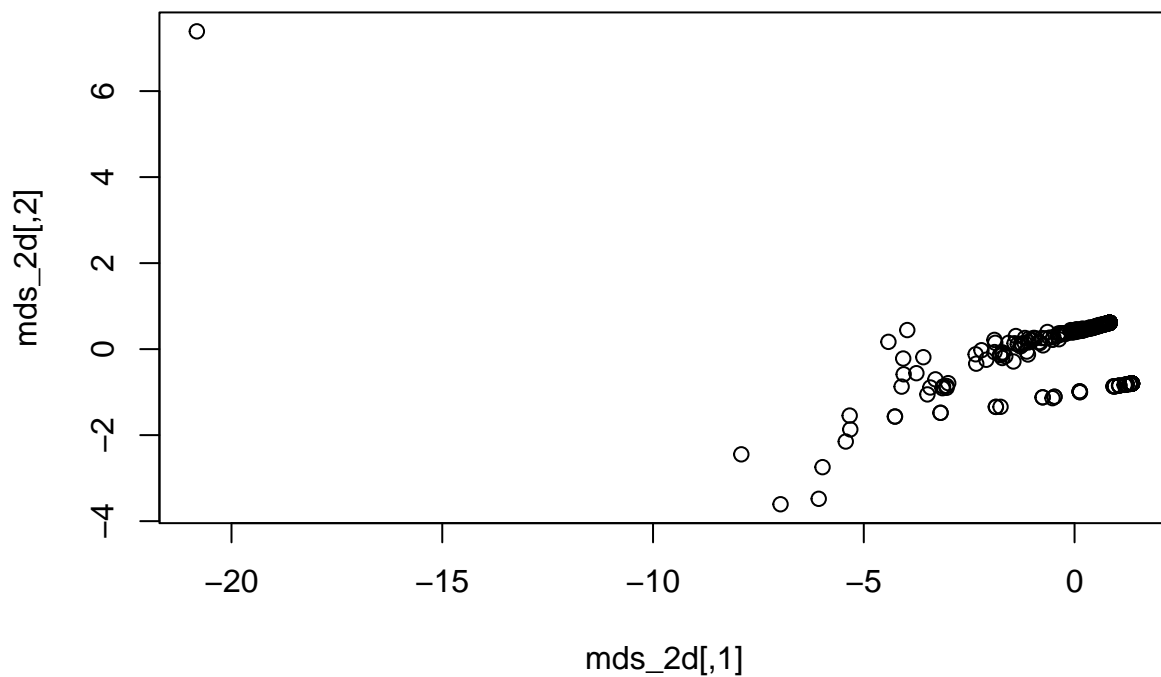
For the second linear regression model of Moderate Dyskaryosis vs Borderline Changes, the coefficient of the slope of the line was significant with an alpha level of  $<0.001$ . The intercept of the line was not significant at any level. The adjusted R-squared value is 0.7823, indicating that 78.23% of the variance in the data can be explained by the model. The F statistic is large with a p-value  $< 0.05$ , indicating that this regression is significant.

For the third linear regression model of Severe Dyskaryosis vs Borderline Changes, the coefficient of the slope of the line was significant with an alpha level of  $<0.001$ . The intercept of the line was not significant at any level. The adjusted R-squared value is 0.8387, indicating that 83.87% of the variance in the data can be explained by the model. The F statistic is large with a p-value of  $< 0.05$ , indicating that this regression is significant.

The regression model of Severe Dyskaryosis vs Borderline Changes appears to be the best at predicting incidence of severe dyskaryosis from borderline changes due to it having the largest F-statistic value (1244) as well as the largest adjusted R-squared value (0.8387). This could be because borderline changes in a neoplasma will typically be larger and more easily detectable in later stage cancers, meaning that large borderline changes would have a higher correlation with severe dyskaryosis.

## Clustering

```
pairwise_dist <- dist(scaled_data)
mds_2d <- cmdscale(pairwise_dist, k = 2)
plot(mds_2d)
```



### BIC MIC Models Raynah Cheng

```
#splitting data into test and train
train_index <- 1:165
test_index <- 166:nrow(comparison_data)

train_frame <- comparison_data[train_index,]
test_frame <- comparison_data[test_index,]

train_model <- lm(Borderline_changes ~ ., data = train_frame)
summary(train_model)
```

```
##
## Call:
## lm(formula = Borderline_changes ~ ., data = train_frame)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5767.6  -992.5  -897.0   595.1 12859.7
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    988.5484    239.2512   4.132 5.77e-05 ***
## Mild_dyskaryosis     0.6606     0.1458   4.531 1.14e-05 ***
## Moderate_dyskaryosis -3.4393     0.8508  -4.042 8.18e-05 ***
```

```
## Severe_dyskaryosis      3.3890      0.3874      8.748 2.77e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2377 on 161 degrees of freedom
## Multiple R-squared:  0.8112, Adjusted R-squared:  0.8076
## F-statistic: 230.5 on 3 and 161 DF,  p-value: < 2.2e-16
```

```
back_BIC <- step(train_model, direction = "backward",
                  k = log(nrow(train_frame)), trace = 0)
summary(back_BIC)
```

```
##
## Call:
## lm(formula = Borderline_changes ~ Mild_dyskaryosis + Moderate_dyskaryosis +
##     Severe_dyskaryosis, data = train_frame)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5767.6  -992.5  -897.0   595.1 12859.7
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    988.5484    239.2512   4.132 5.77e-05 ***
## Mild_dyskaryosis     0.6606     0.1458   4.531 1.14e-05 ***
## Moderate_dyskaryosis -3.4393     0.8508  -4.042 8.18e-05 ***
## Severe_dyskaryosis     3.3890     0.3874   8.748 2.77e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2377 on 161 degrees of freedom
## Multiple R-squared:  0.8112, Adjusted R-squared:  0.8076
## F-statistic: 230.5 on 3 and 161 DF,  p-value: < 2.2e-16
```

```
back_MIC <- step(train_model, direction = "backward",
                  k = nrow(train_frame), trace = 0)
summary(back_MIC)
```

```
##
## Call:
## lm(formula = Borderline_changes ~ Severe_dyskaryosis, data = train_frame)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5789.6 -1266.2  -940.9   535.1  9728.3
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    1275.1192    233.8459   5.453 1.81e-07 ***
## Severe_dyskaryosis     2.6817     0.1093  24.534 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Residual standard error: 2510 on 163 degrees of freedom
## Multiple R-squared:  0.7869, Adjusted R-squared:  0.7856
## F-statistic: 601.9 on 1 and 163 DF,  p-value: < 2.2e-16
```

The BIC model chooses all of the dyskaryosis variables, with equal about equal levels of significance. However, we can see that the severe dyskaryosis t value is higher, at 8.748, compared to the others at around 4.5. For our BIC model, 81.12% of the variance can be explained by our model. The MIC model only chooses the severe dyskaryosis variable, at a .001 level of significance, with a much higher t value than the BIC model. However, the variance that can be explained goes down a little bit, at 78.69%.

```
test_MIC <- predict(back_MIC, test_frame)
test_BIC <- predict(back_BIC, test_frame)

cor(test_MIC, test_frame$Borderline_changes, use = "complete.obs")^2 #r-squared for back_MIC model
```

```
## [1] 0.9967432
```

```
cor(test_BIC, test_frame$Borderline_changes, use = "complete.obs")^2 #r-squared for back_BIC model
```

```
## [1] 0.9911532
```

```
errors_MIC <- test_frame$Borderline_changes - test_MIC
errors_BIC <- test_frame$Borderline_changes - test_BIC

sqrt(mean(errors_MIC^2, na.rm = TRUE)) #RMSE for back_MIC
```

```
## [1] 1763.824
```

```
sqrt(mean(errors_BIC^2, na.rm = TRUE)) #RMSE for back_BIC
```

```
## [1] 1650.034
```

```
mean(abs(errors_MIC), na.rm = TRUE) #MAE for back_MIC
```

```
## [1] 1286.169
```

```
mean(abs(errors_BIC), na.rm = TRUE) #MAE for back_BIC
```

```
## [1] 1056.618
```

Model	$R^2$	RMSE	MAE
MIC	.997	1763.82	1286.17
BIC	.991	1650.03	1056.62

The BIC model has a lower  $R^2$  value, but the errors are lower. The MIC model has a higher  $R^2$  value, but the errors are higher. I think that we would pick the BIC model because the  $R^2$  is not that much lower and it predicts our data a lot better.