# Ordinal logistic model on large, classified windows data

Ruth Gómez Graciani

## Prepare the data

First, we obtain the density distribution, and local minima and maxima for the recombination map.
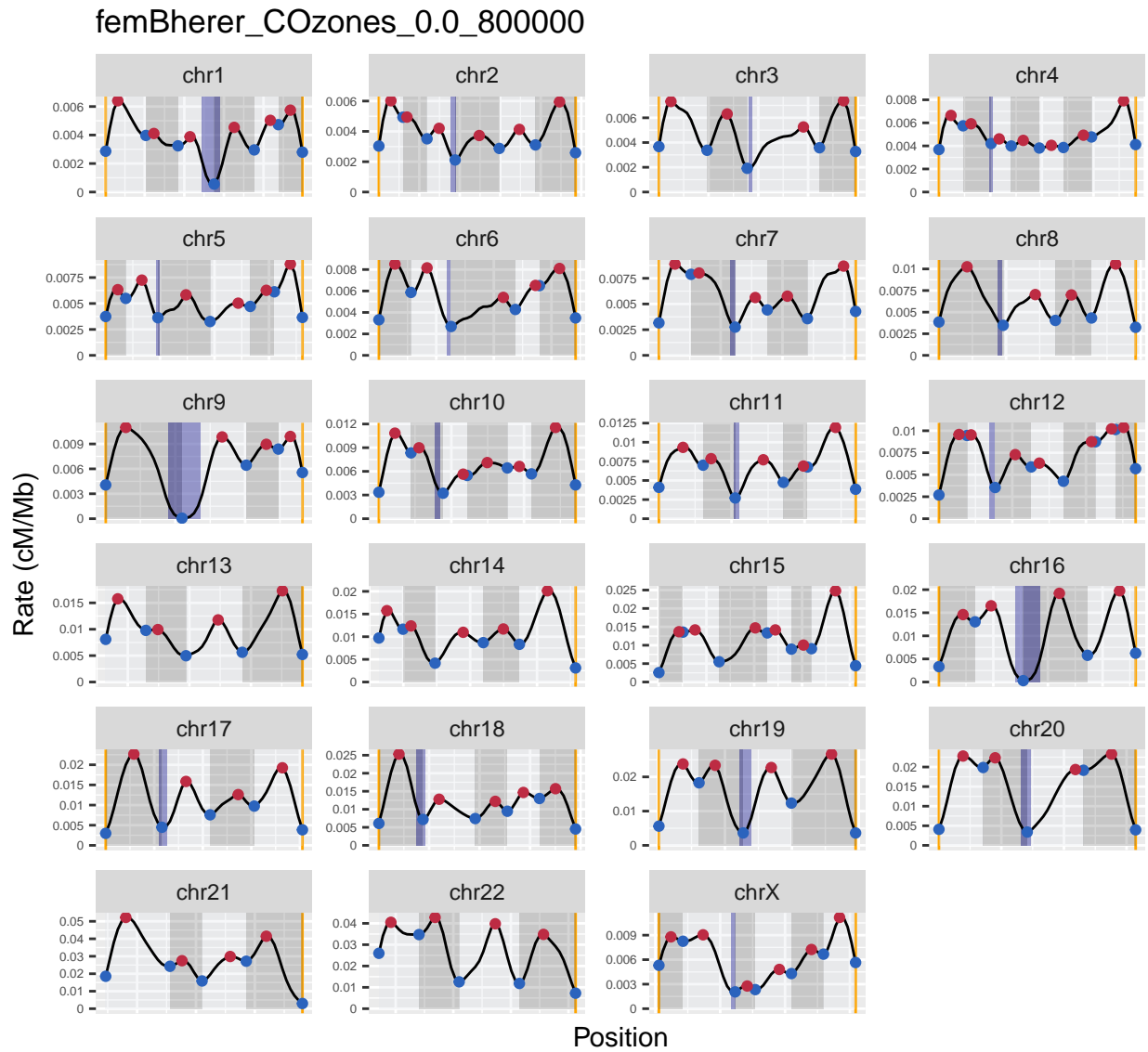


Figure 1: Crossover zones; centromeres in blue, workspace limits in orange.

Next, we define telomeric regions as the space between the chromosome start to the next local minimum, or between the chromosome end to the previous local minimum. We also define centromeric regions as the space between two local maxima that contains the centromere. When the local maximum delimiting a centromeric region is the same as the peak from the corresponding telomeric region (see chr1, chr5, chr7, chr8, etc.), the limit between the telomeric and centromeric regions is defined as the center point between the local maximum corresponding to the telomeric peak and the local minimum corresponding to the centromere valley. These categories will be represented as the "Color" variable in this analysis.
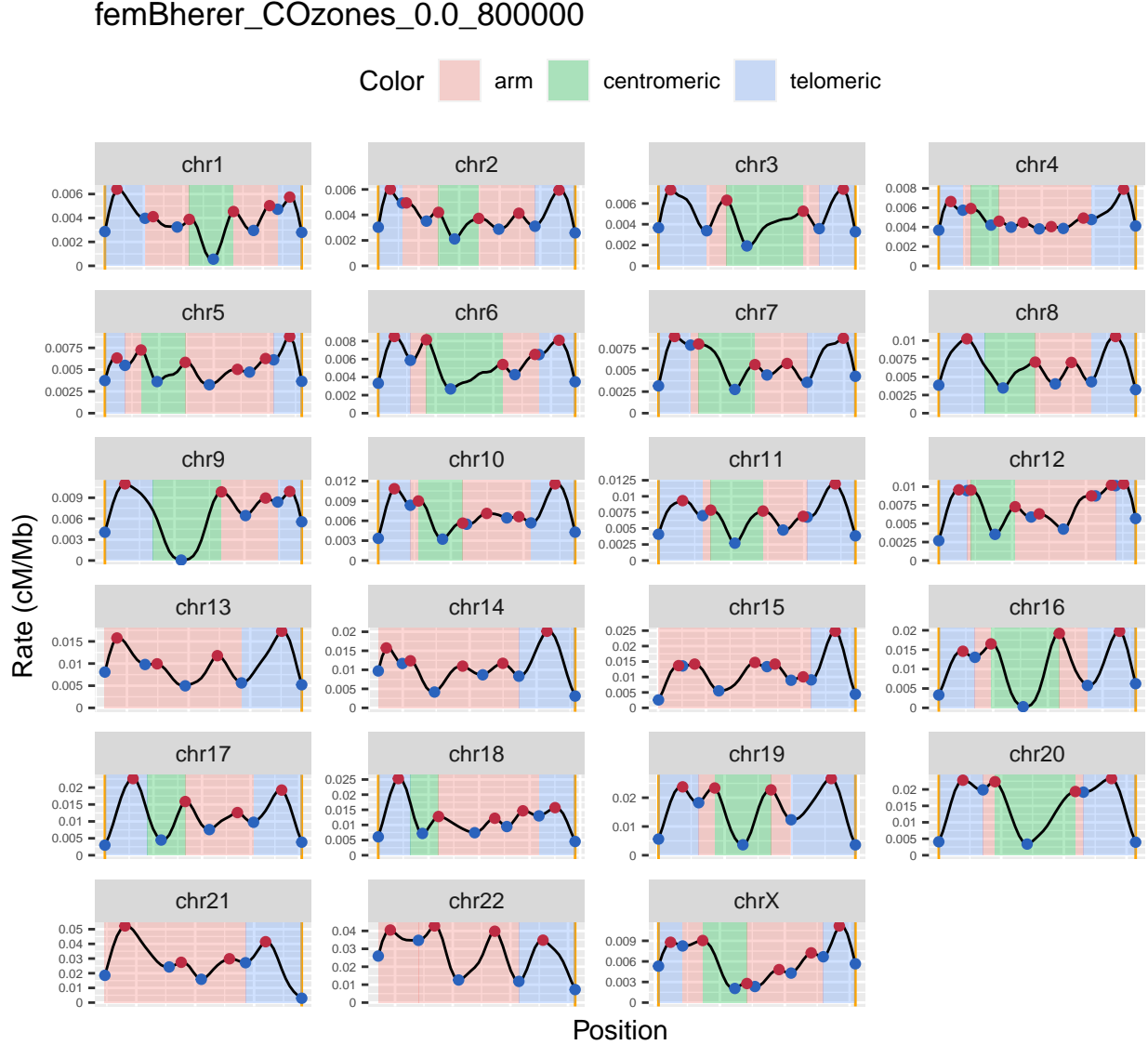


Figure 2: Color-coded windows for telomeric, centromeric and arm categories.

## Descriptive statistics

Raw data:

```
##    Chromosome     Start       End       Color invCenters NHCenters NAHRCenters
## 1      chr10    158946  22251109    telomeric          3         2           1
## 2      chr10  22251109  27774149          arm          0         0           0
## 3      chr10  27774149  58150873  centromeric          2         1           1
## 4      chr10  58150873 105096718          arm          3         3           0
## 5      chr10 105096718 135473442    telomeric          1         1           0
## 6      chr11    241489  30481001    telomeric          2         1           1
##    Length.Mb allRepCounts WAvgRate.perMb
## 1 22.092163          340       1.855788
## 2  5.523041          200       1.465803
## 3 30.376724         2350       1.215003
## 4 46.945846          808       1.306162
## 5 30.376724          200       1.864682
## 6 30.239512          748       1.591677
```

For each window, I calculated the number of total inversions, NH inversions, and NAHR inversions, the window length in Mb, number of repeats and the average recombination rate in cM/Mb.

I want to perform Ordinal Logistic Regressions on different subsets of the data. The assumptions of the Ordinal Logistic Regression are as follow:

1. The dependent variable is ordered.
2. One or more of the independent variables are either continuous, categorical or ordinal.
3. No multi-collinearity.
4. Proportional odds.

I show the data distributions in the figure below. The inversion counts have only a number of possible options, so they can be considered an ordinal variable. The independent variables are continuous and categorical, so assumptions 1 and 2 are satisfied
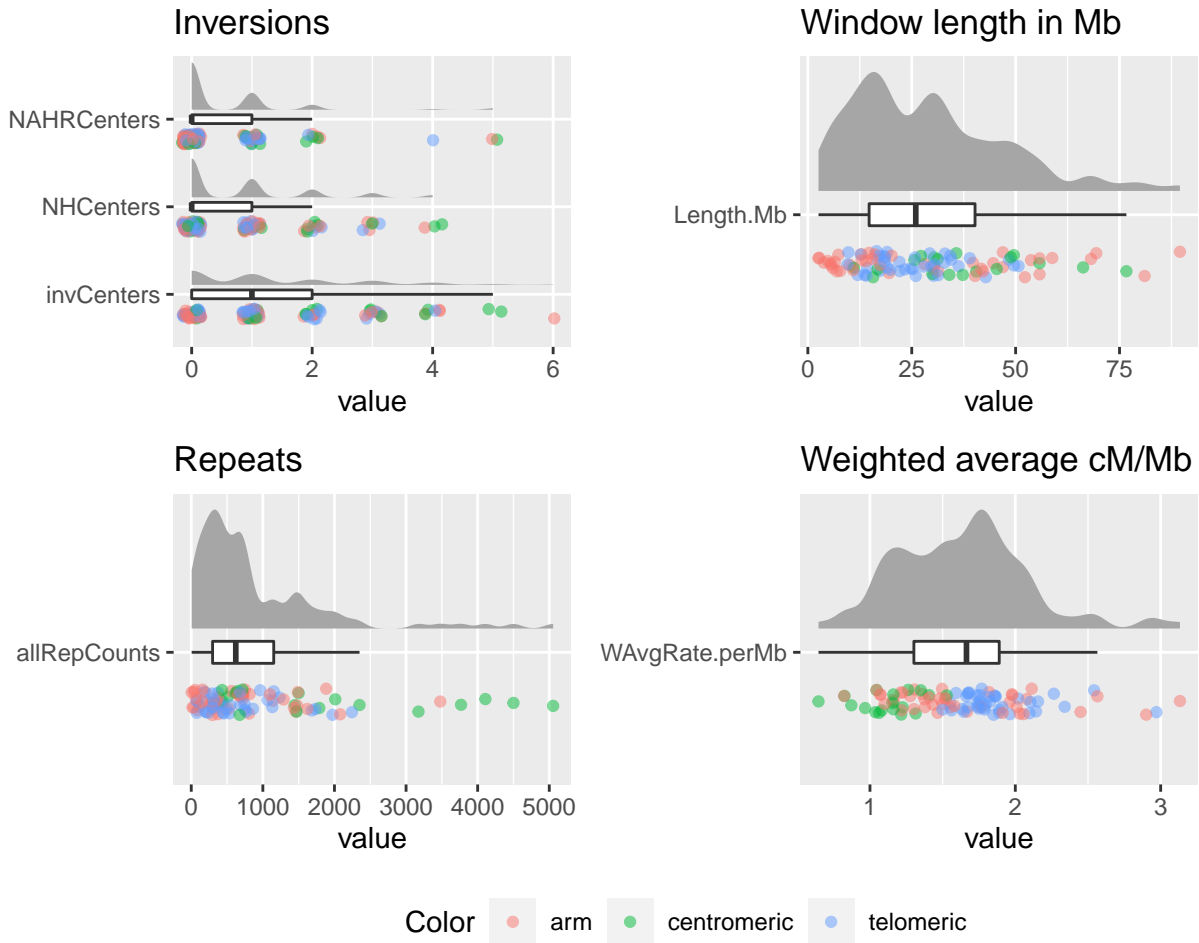
Figure 3: Distribution of variables.

We see that some categories have low number of cases, so I will make a "3 or more" category when relevant.

```
## [1] "Original counts"

##   CountGroups invCenters NHCenters NAHRCenters
## 1           0         38        55          67
## 2           1         29        26          24
## 3           2         14        11           7
## 4           3         11         6          NA
## 5           4          6         3           1
## 6           5          2        NA           2
## 7           6          1        NA          NA

## [1] "New counts"

##   CountGroups invCategory NHCategory NAHRCategory
## 1           0          38         55           67
## 2           1          29         26           24
## 3           2          14         11            7
## 4          3+          20          9            3
```

With these groups, I visualize the relationships between dependent and independent variables.

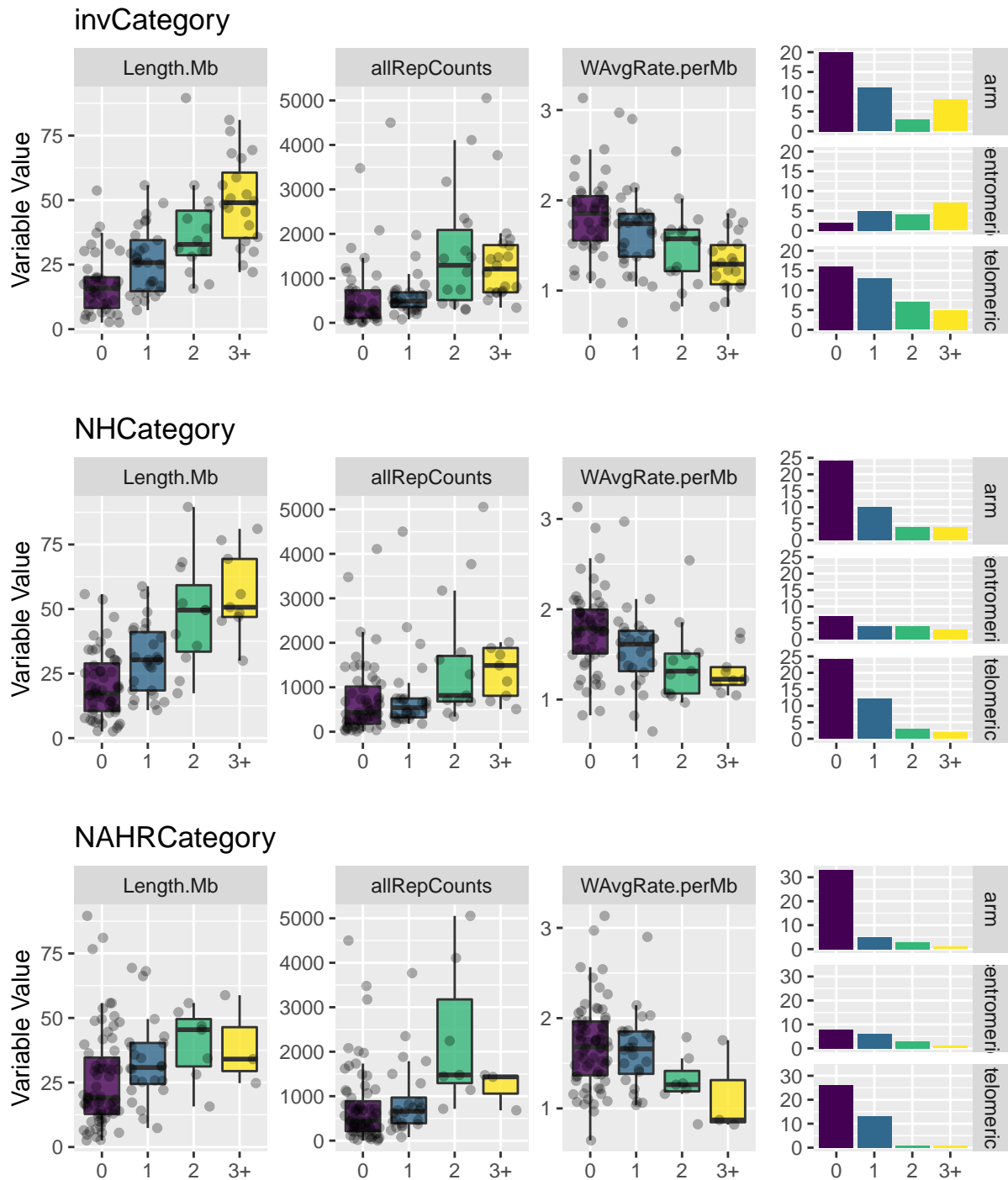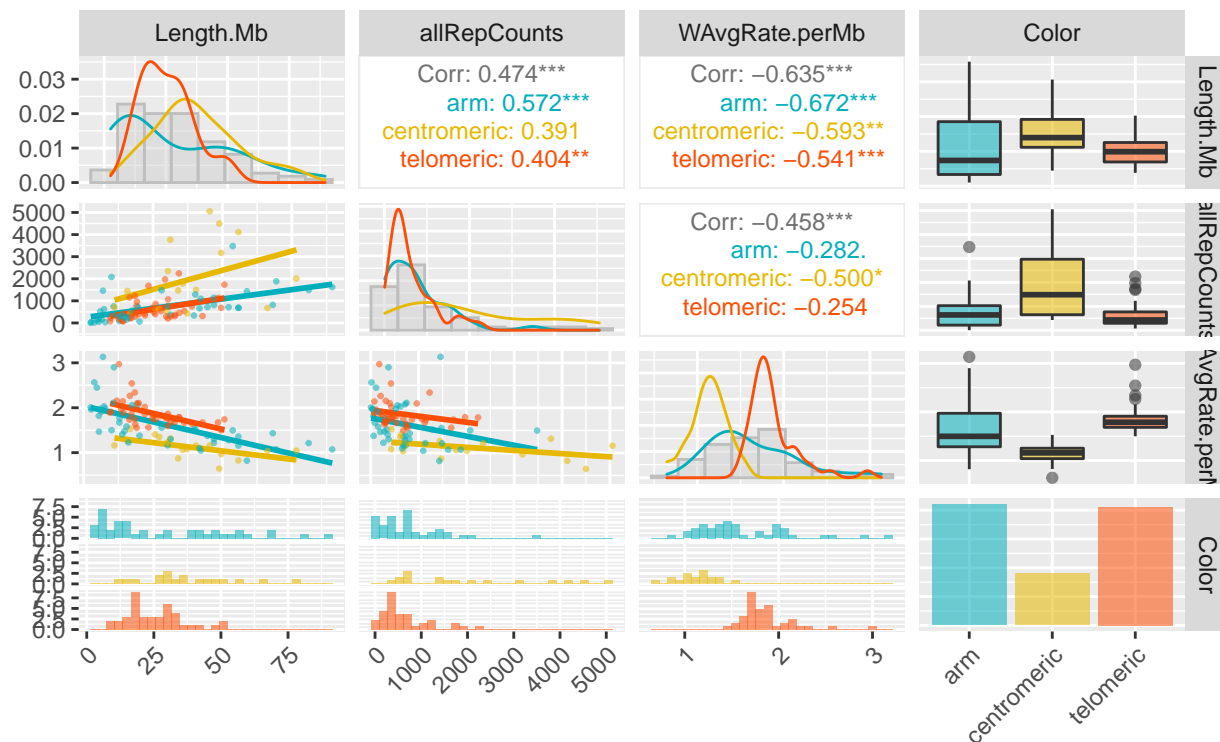# Differences in each chromosomal variable between inversion count groups



Figure 4: Potential effect of independent variables on the different types of invesions.

Finally, I will test assumption number 3, no multi-collinearity between independent variables.
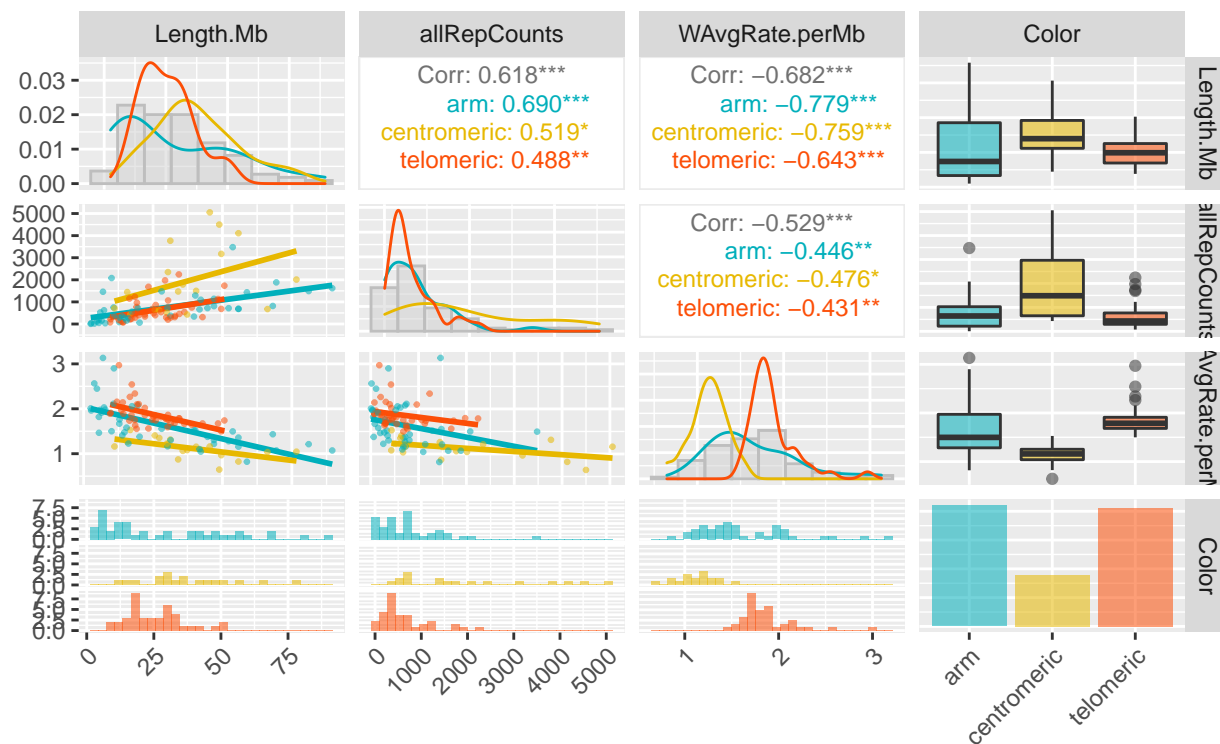
## Pearson correlation



## Spearman correlation



Figure 5: Correlations between variables.

We see that our three variables are significantly correlated, but this does not confirm multi-collinearity. I perform a variance inflation factor test on the corresponging linear model to further check the multi-collinearity.

```
##                  GVIF Df GVIF^(1/(2*Df))
## Length.Mb      2.001170  1        1.414627
## allRepCounts  1.647309  1        1.283475
## Color          1.863293  2        1.168343
## WAvgRate.perMb 2.446559  1        1.564148
```

The general rule of thumbs for VIF test is that if the VIF value is greater than 10, then there is multi-collinearity, so we can say that the third assumption (no multi-collinearity) is satisfied.

The proportional odds assumption will be tested for each model that we fit in the following analyses.

## Variable scalation (optional)

Standardized coefficients are useful in our case to compare effects of predictors reported in different units. The most straightforward way is using the Agresti method of standardization, applied with the `scale()` function.

```
##     Length.Mb      Length.Mb.Scaled   allRepCounts     allRepCounts.Scaled
## Min.    : 2.565   Min.    :-1.4230   Min.    :   6.0   Min.    :-0.9215
## 1st Qu.:14.738   1st Qu.:-0.7697   1st Qu.: 298.0   1st Qu.:-0.6226
## Median :25.971   Median :-0.1668   Median : 622.0   Median :-0.2911
## Mean   :29.078   Mean    : 0.0000   Mean    : 906.4   Mean    : 0.0000
## 3rd Qu.:40.233   3rd Qu.: 0.5987   3rd Qu.:1152.0   3rd Qu.: 0.2513
## Max.   :89.545   Max.    : 3.2453   Max.    :5054.0   Max.    : 4.2445
## WAvgRate.perMb  WAvgRate.perMb.Scaled
## Min.   :0.646   Min.    :-2.15076
## 1st Qu.:1.302   1st Qu.:-0.72706
## Median :1.664   Median : 0.05816
## Mean   :1.637   Mean    : 0.00000
## 3rd Qu.:1.889   3rd Qu.: 0.54701
## Max.   :3.133   Max.    : 3.24489
```

Once the model is fitted, we can use the sd to transform scaled coefficients to natural coefficients and viceversa.

## Total inversions (invCategory)

**Model fitting**

```
## Call:
## polr(formula = myFormula, data = winRegions, Hess = T)
##
## Coefficients:
##                      Value Std. Error t value
## Length.Mb         9.952e-02  0.0142331  6.9919
## allRepCounts     -8.528e-05  0.0003422 -0.2492
## Colorcentromeric  1.651e+00  0.6350390  2.5998
## Colortelomeric    3.211e-01  0.4621403  0.6948
## WAvgRate.perMb    8.173e-01  0.2038465  4.0096
##
## Intercepts:
##      Value   Std. Error t value
## 0|1  3.6692  0.1448     25.3385
## 1|2  5.4864  0.3268     16.7867
## 2|3+ 6.6198  0.4170     15.8734
##
## Residual Deviance: 208.6729
## AIC: 224.6729
```

We compare the t-value against the standard normal distribution to calculate the p-value.

```
##                        Value    Std. Error    t value     p value
## Length.Mb          9.951611e-02 0.0142330757  6.9918906 0.00000000
## allRepCounts      -8.527591e-05 0.0003422055 -0.2491950 0.80320991
## Colorcentromeric   1.650947e+00 0.6350390465  2.5997565 0.00932899
## Colortelomeric     3.211146e-01 0.4621403074  0.6948422 0.48715418
## WAvgRate.perMb     8.173350e-01 0.2038465195  4.0095609 0.00006083
## 0|1                3.669234e+00 0.1448086001 25.3385087 0.00000000
## 1|2                5.486432e+00 0.3268325475 16.7866749 0.00000000
## 2|3+               6.619768e+00 0.4170341732 15.8734428 0.00000000
```

We can also get confidence intervals for the parameter estimates. These can be obtained either by profiling the likelihood function or by using the standard errors and assuming a normal distribution. Note that profiled CIs are not symmetric (although they are usually close to symmetric). If the 95% CI does not cross 0, the parameter estimate is statistically significant.

```
## [1] "Profiling likelihod"
```

```
##                        2.5 %         97.5 %
## Length.Mb          0.0625867701 0.1413139203
## allRepCounts      -0.0006388095 0.0004698302
## Colorcentromeric   0.2508033684 3.1302149241
## Colortelomeric    -0.5861609138 1.2513784295
## WAvgRate.perMb             NA            NA
```

```
## [1] "Assuming a normal distribtuion"
```

```
##                        2.5 %         97.5 %
## Length.Mb          0.0716197926 0.1274124241
## allRepCounts      -0.0007559864 0.0005854345
## Colorcentromeric   0.4062932084 2.8956005283
## Colortelomeric    -0.5846637590 1.2268929576
```

```
## WAvgRate.perMb    0.4178031930 1.2168668663
```

We convert the coefficients into odds ratios. To get the OR and confidence intervals, we just exponentiate the estimates and confidence intervals (here I used the likelihood confidence intervals).

```
##                   Odds Ratio      2.5%      97.5%
## Length.Mb          1.1046363 1.0645868  1.151786
## allRepCounts       0.9999147 0.9993614  1.000470
## Colorcentromeric   5.2119125 1.2850574 22.878896
## Colortelomeric     1.3786636 0.5564595  3.495157
## WAvgRate.perMb     2.2644571        NA        NA
```

Example of interpretation: "For 1 unit increase in Length.Mb, a window is 1.1046363 times more likely to increase in inversion amount category."
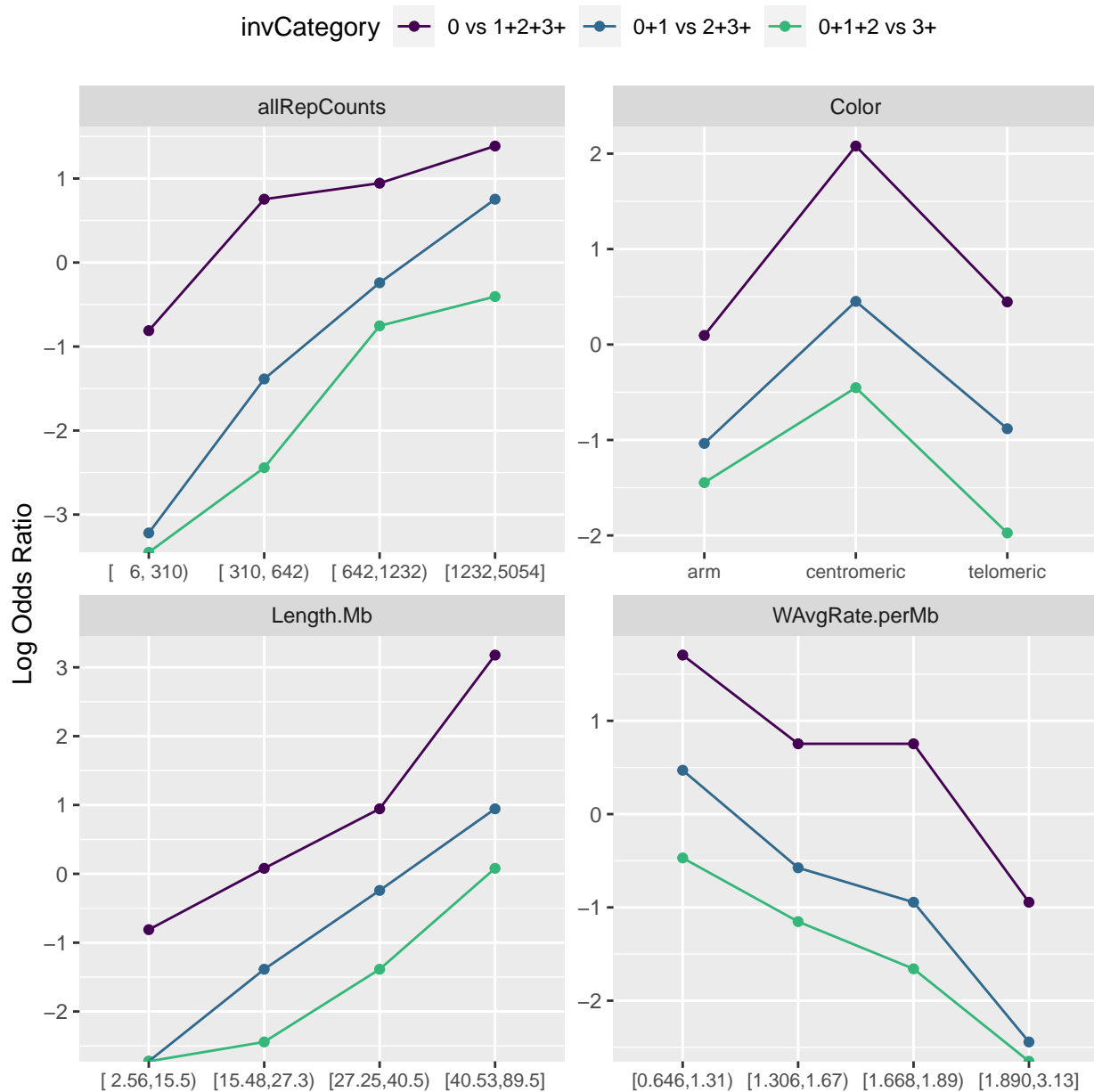
**Proportional odds assessment**

Now we should test the proportional odds or parallel regression assumption. If it is satisfied, the coefficients are valid for all the cases (i.e. the same coefficient is valid for increasing from 0 to 1 inversions, from 1 to 2, etc.). If this assumption is violated, different models are needed to describe the relationship between each pair of outcome groups.

We test the parallel regression assumption with a Brant test:

```
## -------------------------------------------------------
## Test for      X2  df  probability
## -------------------------------------------------------
## Omnibus          4.66   10  0.91
## Length.Mb        0.68    2  0.71
## allRepCounts     1.55    2  0.46
## Colorcentromeric 0.39    2  0.82
## Colortelomeric   1.37    2   0.51
## WAvgRate.perMb   0.26    2   0.88
## -------------------------------------------------------
##
## H0: Parallel Regression Assumption holds
```

We can also evaluate the parallel regression visually. We transform the ordinal dependent variable with k categories into a series of k-1 binary variables that indicate whether the dependent value is above or below a cutpoint (e.g. windows with at least 2 inversions vs windows with less than 2 inversions). We then calculate the observed Log Odds Ratio for each binary variable across multiple value ranges of the independent variables. The lines should be approximately parallel, that each independent variable affects the probability of increasing by 1 level the inversion count in the same way, for all transitions, and that we don't need a specific model for each level increase.

## Proportional odds visual test

## Predicted probabilites

Although our objective is to describe the dataset, predicted probabilities are usually easier to understand than either the coefficients or the Odds Ratios.
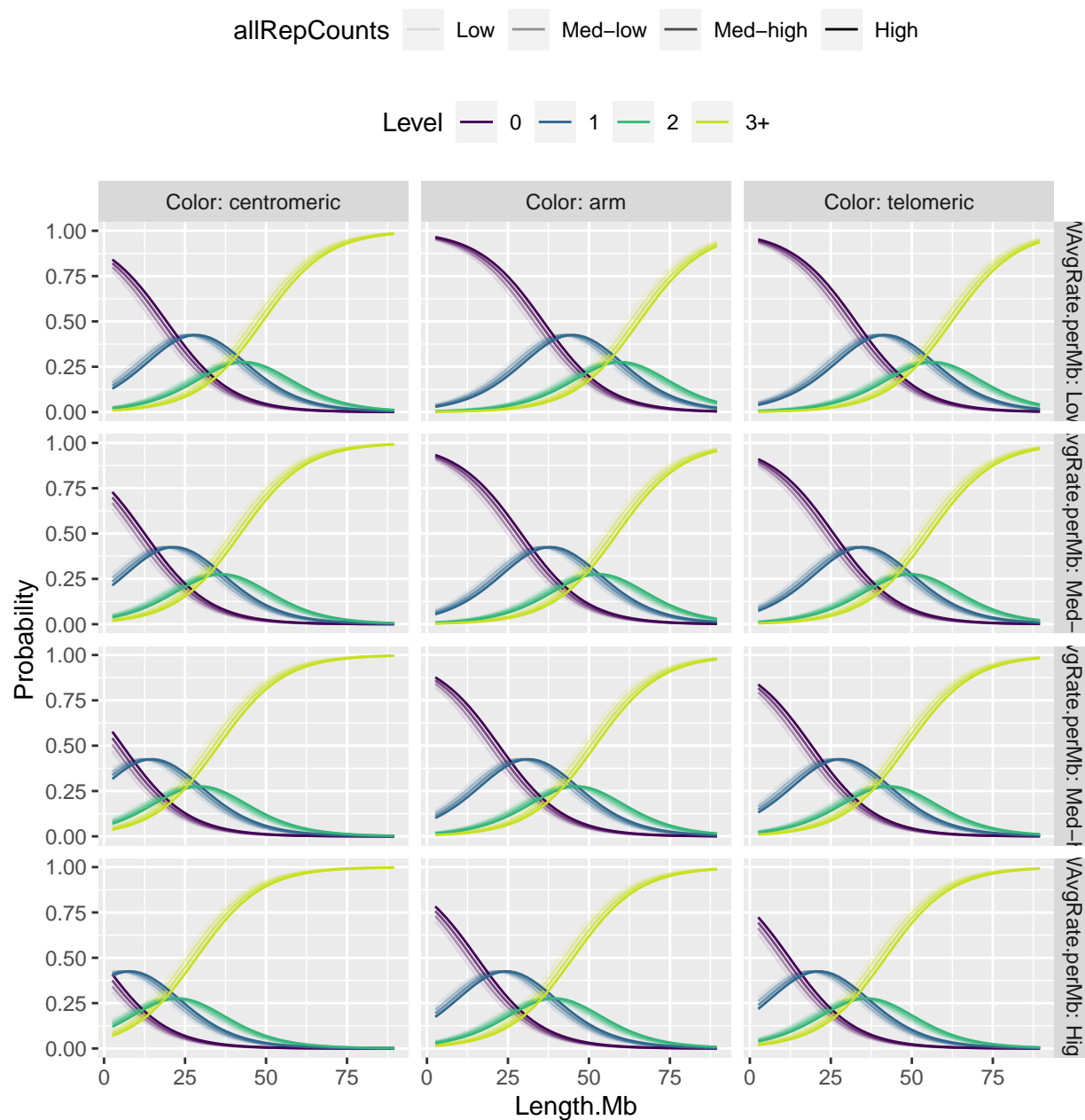


Figure 6: Probabiilty of having 0 to >3 inversions depending on multiple independent variables

## Total inversions (NHCategory)

**Model fitting**

```
## Call:
## polr(formula = myFormula, data = winRegions, Hess = T)
##
## Coefficients:
##                      Value Std. Error  t value
## Length.Mb          9.143e-02  0.0124309  7.35470
## allRepCounts      -3.228e-05  0.0003282 -0.09835
## Colorcentromeric   9.362e-01  0.6286970  1.48906
## Colortelomeric     4.652e-02  0.5001710  0.09300
## WAvgRate.perMb     8.872e-01  0.2270231  3.90789
##
## Intercepts:
##       Value   Std. Error t value
## 0|1   4.3567  0.1413     30.8398
## 1|2   6.1522  0.3573     17.2184
## 2|3+  7.5650  0.5164     14.6505
##
## Residual Deviance: 182.8368
## AIC: 198.8368
```

We compare the t-value against the standard normal distribution to calculate the p-value.

```
##                       Value       Std. Error      t value      p value
## Length.Mb          9.142563e-02 0.0124309048   7.35470454 0.00000000
## allRepCounts      -3.228081e-05 0.0003282103  -0.09835403 0.92165118
## Colorcentromeric   9.361700e-01 0.6286969984   1.48906384 0.13647056
## Colortelomeric     4.651670e-02 0.5001709948   0.09300159 0.92590230
## WAvgRate.perMb     8.871811e-01 0.2270230851   3.90788934 0.00009311
## 0|1                4.356732e+00 0.1412699284  30.83977137 0.00000000
## 1|2                6.152218e+00 0.3573051852  17.21838458 0.00000000
## 2|3+               7.565027e+00 0.5163680997  14.65045312 0.00000000
```

We can also get confidence intervals for the parameter estimates. These can be obtained either by profiling the likelihood function or by using the standard errors and assuming a normal distribution. Note that profiled CIs are not symmetric (although they are usually close to symmetric). If the 95% CI does not cross 0, the parameter estimate is statistically significant.

```
## [1] "Profiling likelihod"
```

```
##                        2.5 %         97.5 %
## Length.Mb          0.0572532547 0.1298089217
## allRepCounts      -0.0005916123 0.0005087736
## Colorcentromeric  -0.4550261982 2.3681680663
## Colortelomeric    -0.9293705517 1.0588746546
## WAvgRate.perMb              NA           NA
```

```
## [1] "Assuming a normal distribtuion"
```

```
##                        2.5 %         97.5 %
## Length.Mb          0.0670615065 0.1157897580
## allRepCounts      -0.0006755612 0.0006109996
## Colorcentromeric  -0.2960535098 2.1683934383
## Colortelomeric    -0.9338004378 1.0268338340
```

```
## WAvgRate.perMb     0.4422240241 1.3321381650
```

We convert the coefficients into odds ratios. To get the OR and confidence intervals, we just exponentiate the estimates and confidence intervals (here I used the likelihood confidence intervals).

```
##                    Odds Ratio      2.5%      97.5%
## Length.Mb          1.0957353 1.0589240  1.138611
## allRepCounts       0.9999677 0.9994086  1.000509
## Colorcentromeric   2.5501954 0.6344313 10.677813
## Colortelomeric     1.0476156 0.3948021  2.883125
## WAvgRate.perMb     2.4282749        NA        NA
```

Example of interpretation: "For 1 unit increase in Length.Mb, a window is 1.0957353 times more likely to increase in inversion amount category."
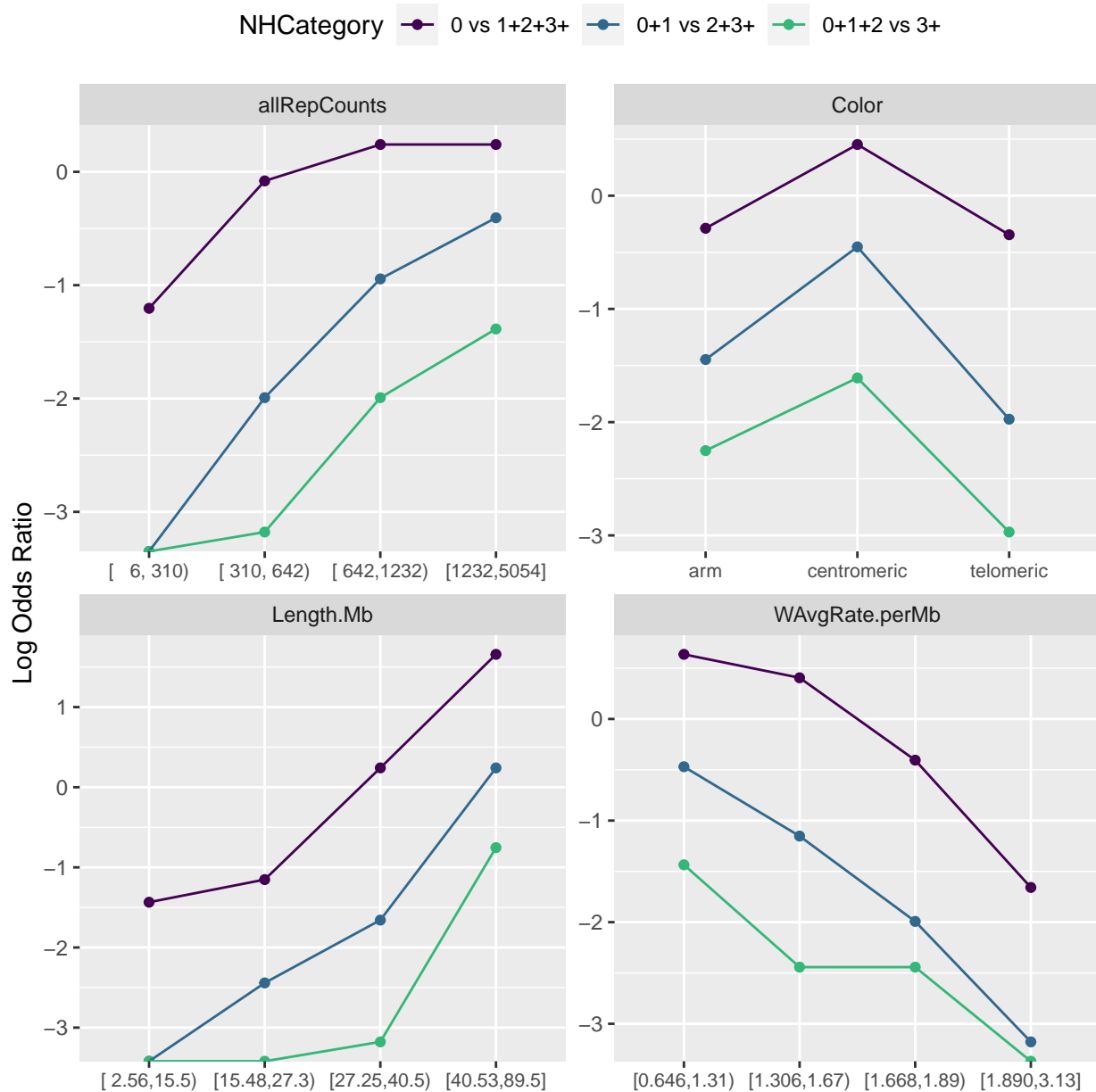
### Proportional odds assessment

Now we should test the proportional odds or parallel regression assumption. If it is satisfied, the coefficients are valid for all the cases (i.e. the same coefficient is valid for increasing from 0 to 1 inversions, from 1 to 2, etc.). If this assumption is violated, different models are needed to describe the relationship between each pair of outcome groups.

We test the parallel regression assumption with a Brant test:

```
## --------------------------------------------------------
## Test for      X2  df  probability
## --------------------------------------------------------
## Omnibus          11.63   10  0.31
## Length.Mb        5.2 2   0.07
## allRepCounts     1.93    2   0.38
## Colorcentromeric 4.56    2   0.1
## Colortelomeric       0.04    2   0.98
## WAvgRate.perMb       8.22    2   0.02
## --------------------------------------------------------
##
## H0: Parallel Regression Assumption holds
```

We can also evaluate the parallel regression visually. We transform the ordinal dependent variable with k categories into a series of k-1 binary variables that indicate whether the dependent value is above or below a cutpoint (e.g. windows with at least 2 inversions vs windows with less than 2 inversions). We then calculate the observed Log Odds Ratio for each binary variable across multiple value ranges of the independent variables. The lines should be approximately parallel, that each independent variable affects the probability of increasing by 1 level the inversion count in the same way, for all transitions, and that we don't need a specific model for each level increase.

## Proportional odds visual test

**Predicted probabilites**

Although our objective is to describe the dataset, predicted probabilities are usually easier to understand than either the coefficients or the Odds Ratios.
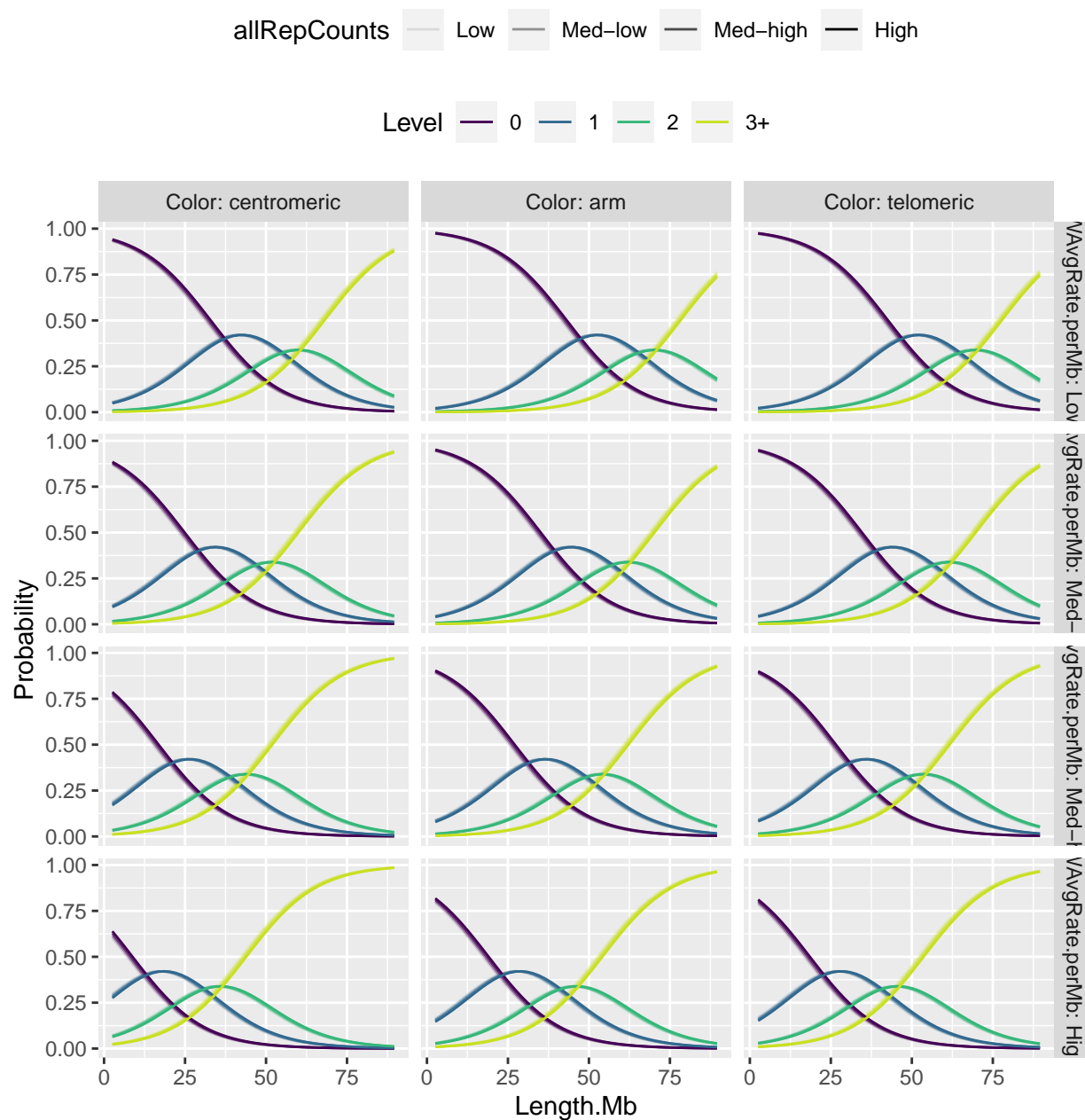


Figure 7: Probabiilty of having 0 to >3 inversions depending on multiple independent variables

## Total inversions (NAHRCategory)

**Model fitting**

```
## Call:
## polr(formula = myFormula, data = winRegions, Hess = T)
##
## Coefficients:
##                       Value Std. Error t value
## Length.Mb          0.0188601   0.011921   1.582
## allRepCounts       0.0003143   0.000348   0.903
## Colorcentromeric   0.9209545   0.630055   1.462
## Colortelomeric     0.8544856   0.541881   1.577
## WAvgRate.perMb    -0.2531833   0.245667  -1.031
##
## Intercepts:
##        Value   Std. Error t value
## 0|1    1.6899  0.1483     11.3962
## 1|2    3.3881  0.3617      9.3667
## 2|3+   4.7416  0.6271      7.5614
##
## Residual Deviance: 168.9569
## AIC: 184.9569
```

We compare the t-value against the standard normal distribution to calculate the p-value.

```
##                        Value    Std. Error    t value    p value
## Length.Mb         0.0188600857 0.0119213586  1.5820416 0.1136401
## allRepCounts      0.0003142673 0.0003480284  0.9029934 0.3665294
## Colorcentromeric  0.9209544998 0.6300554837  1.4617038 0.1438224
## Colortelomeric    0.8544856015 0.5418806717  1.5768889 0.1148211
## WAvgRate.perMb   -0.2531832611 0.2456665356 -1.0305973 0.3027297
## 0|1               1.6898869893 0.1482844972 11.3962486 0.0000000
## 1|2               3.3881480780 0.3617230218  9.3666918 0.0000000
## 2|3+              4.7416499999 0.6270900119  7.5613547 0.0000000
```

We can also get confidence intervals for the parameter estimates. These can be obtained either by profiling the likelihood function or by using the standard errors and assuming a normal distribution. Note that profiled CIs are not symmetric (although they are usually close to symmetric). If the 95% CI does not cross 0, the parameter estimate is statistically significant.

```
## [1] "Profiling likelihod"
```

```
##                        2.5 %         97.5 %
## Length.Mb         -0.0131523892 0.0520801011
## allRepCounts      -0.0001947444 0.0008165678
## Colorcentromeric  -0.5092010829 2.3735719224
## Colortelomeric    -0.1797715268 1.9844609079
## WAvgRate.perMb              NA           NA
```

```
## [1] "Assuming a normal distribtuion"
```

```
##                        2.5 %         97.5 %
## Length.Mb         -0.0045053477 0.0422255191
## allRepCounts      -0.0003678558 0.0009963905
## Colorcentromeric  -0.3139315567 2.1558405562
## Colortelomeric    -0.2075809990 1.9165522020
```

```
## WAvgRate.perMb   -0.7346808229 0.2283143008
```

We convert the coefficients into odds ratios. To get the OR and confidence intervals, we just exponentiate the estimates and confidence intervals (here I used the likelihood confidence intervals).

```
##                  Odds Ratio      2.5%      97.5%
## Length.Mb         1.0190391 0.9869337   1.053460
## allRepCounts      1.0003143 0.9998053   1.000817
## Colorcentromeric  2.5116867 0.6009755  10.735671
## Colortelomeric    2.3501651 0.8354611   7.275124
## WAvgRate.perMb    0.7763256        NA         NA
```

Example of interpretation: "For 1 unit increase in Length.Mb, a window is 1.0190391 times more likely to increase in inversion amount category."
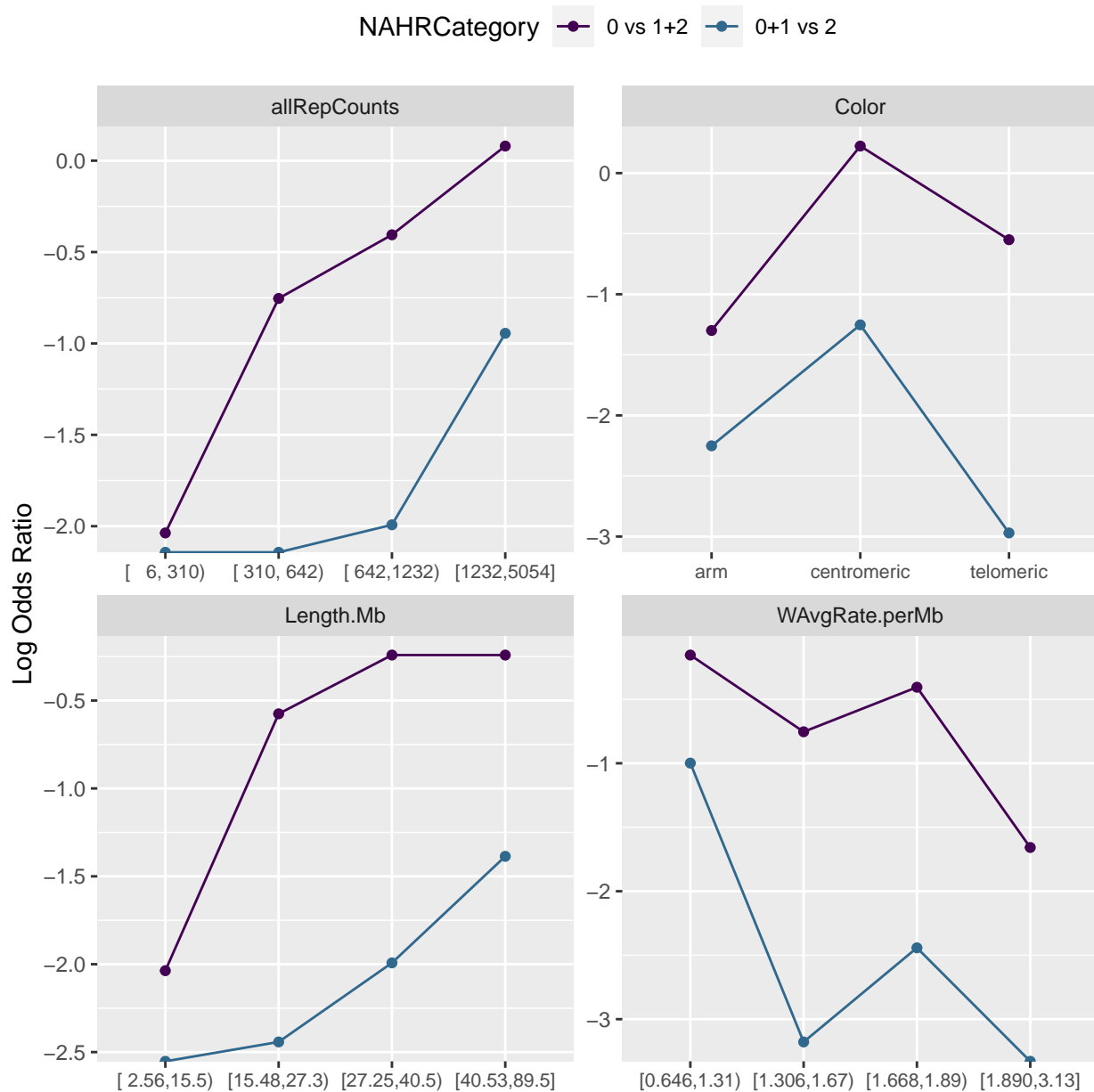
**Proportional odds assessment**

Now we should test the proportional odds or parallel regression assumption. If it is satisfied, the coefficients are valid for all the cases (i.e. the same coefficient is valid for increasing from 0 to 1 inversions, from 1 to 2, etc.). If this assumption is violated, different models are needed to describe the relationship between each pair of outcome groups.

We test the parallel regression assumption with a Brant test:

```
## --------------------------------------------------------
## Test for     X2   df  probability
## --------------------------------------------------------
## Omnibus         15.05   10  0.13
## Length.Mb        3.95    2  0.14
## allRepCounts     7.24    2  0.03
## Colorcentromeric 5.41    2  0.07
## Colortelomeric   3.9 2  0.14
## WAvgRate.perMb   6.12    2   0.05
## --------------------------------------------------------
##
## H0: Parallel Regression Assumption holds
```

We can also evaluate the parallel regression visually. We transform the ordinal dependent variable with k categories into a series of k-1 binary variables that indicate whether the dependent value is above or below a cutpoint (e.g. windows with at least 2 inversions vs windows with less than 2 inversions). We then calculate the observed Log Odds Ratio for each binary variable across multiple value ranges of the independent variables. The lines should be approximately parallel, that each independent variable affects the probability of increasing by 1 level the inversion count in the same way, for all transitions, and that we don't need a specific model for each level increase.



Proportional odds visual test

**Predicted probabilites**

Although our objective is to describe the dataset, predicted probabilities are usually easier to understand than either the coefficients or the Odds Ratios.
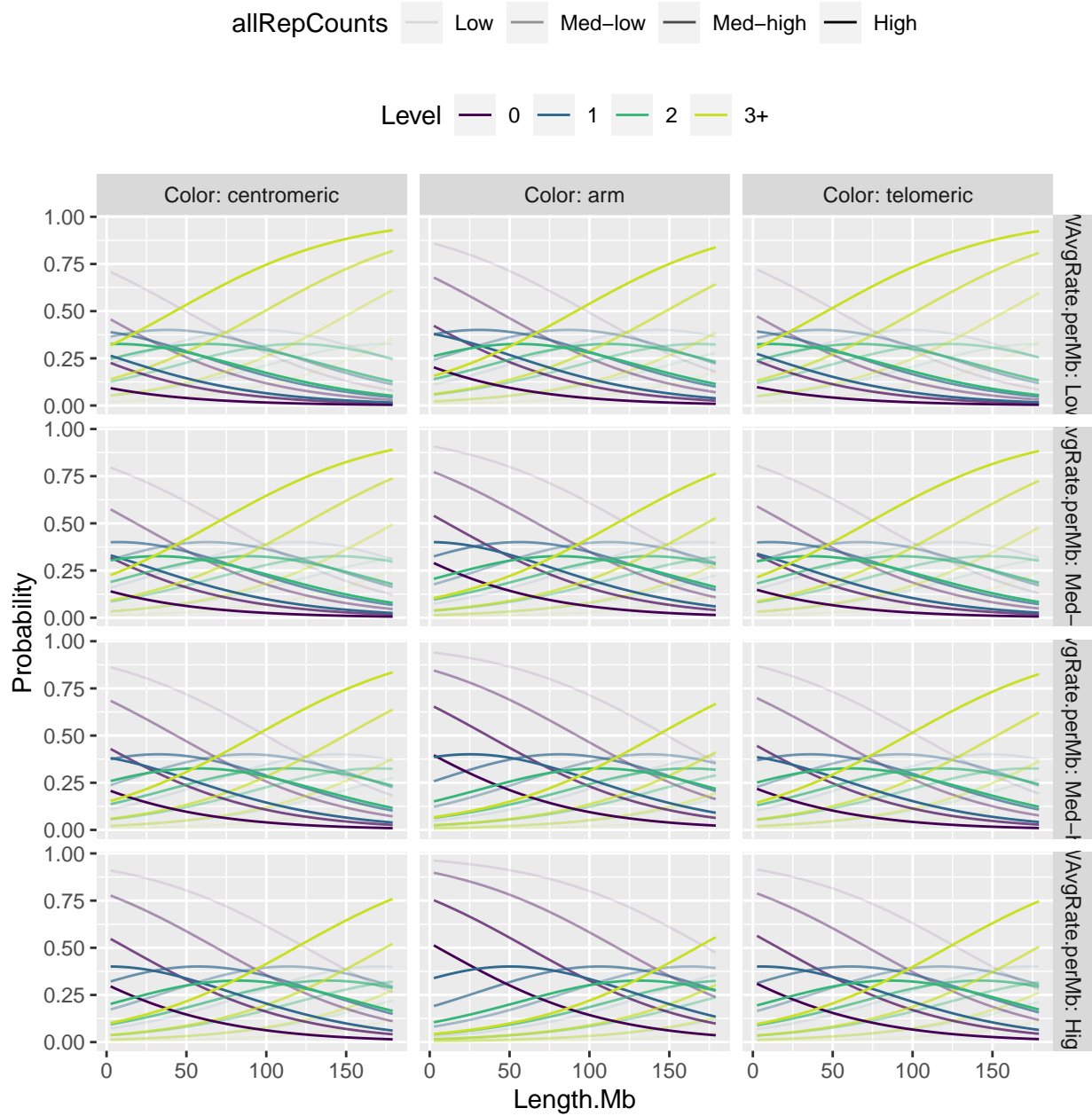


Figure 8: Probabiilty of having 0 to >3 inversions depending on multiple independent variables