**Problem 1**

a.

$$\mathrm{E}[SS_A] = \mathrm{E}\left[rb\sum_{j=1}^{a}(y_{.j.} - y_{...})^2]\right]$$

$$\mathrm{E}[SS_A] = rb\sum_{j=1}^{a}\mathrm{E}[((t_{j.} + e_{.j.}) - (t_{..} + e_{...}))^2]$$

$$\mathrm{E}[SS_A] = rb\sum_{j=1}^{a}\mathrm{E}[(p_j + e_{.j.} - e_{...})^2]$$

$$\mathrm{E}[SS_A] = rb\sum_{j=1}^{a}\left(\mathrm{E}[p_j^2] + 2\mathrm{E}[p_j e_{.j.}] - 2\mathrm{E}[p_j e_{...}] + \mathrm{E}[e_{.j.}^2] - 2\mathrm{E}[e_{.j.}e_{...}] + \mathrm{E}[e_{...}^2]\right)$$

$$\mathrm{E}[SS_A] = rb\sum_{j=1}^{a}\left(p_j^2 + \mathrm{E}\left[\left(\frac{\sum_{i,k}e_{ijk}}{rb}\right)^2\right] - 2\mathrm{E}\left[\frac{\sum_{i,k}e_{ijk}}{rb}\frac{\sum_{i,j,k}e_{ijk}}{rab}\right] + \mathrm{E}\left[\left(\frac{\sum_{i,j,k}e_{ijk}}{rab}\right)^2\right]\right)$$

$$\mathrm{E}[SS_A] = rb\sum_{j=1}^{a}\left(p_j^2 + \frac{\sigma^2}{rb} - \frac{2\sigma^2}{rab} + \frac{\sigma^2}{rab}\right)$$

$$\mathrm{E}[SS_A] = (a-1)\sigma^2 + rb\sum_{j=1}^{a}p_j^2$$

b.

$$\mathrm{E}[SS_{A\times B}] = \mathrm{E}\left[r\sum_{j,k}(y_{.jk} - y_{.j.} - y_{..k} + y_{...})^2]\right]$$

$$\mathrm{E}[SS_{A\times B}] = r\sum_{j,k}\mathrm{E}[((pq)_{jk} + e_{.jk} - e_{.j.} - e_{..k} + e_{...})^2]$$

$$\mathrm{E}[SS_{A\times B}] = r\sum_{j,k}\left[(pq)_{jk} + \mathrm{E}[e_{.jk}^2] + \mathrm{E}[e_{.j.}^2] + \mathrm{E}[e_{..k}^2] + \mathrm{E}[e_{...}^2] - 2\mathrm{E}[e_{.jk}e_{.j.}] - 2\mathrm{E}[e_{.jk}e_{..k}] + 2\mathrm{E}[e_{.jk}e_{...}]\right.$$

$$\left. + 2\mathrm{E}[e_{.j.}e_{..k}] - 2\mathrm{E}[e_{.j.}e_{...}] - 2\mathrm{E}[e_{..k}e_{...}]\right]$$

$$\mathrm{E}[SS_{A\times B}] = r\sum_{j,k}\left[(pq)_{jk} + \frac{\sigma^2}{r} + \frac{\sigma^2}{ra} + \frac{\sigma^2}{rb} + \frac{\sigma^2}{rab} - \frac{2\sigma^2}{rb} - \frac{2\sigma^2}{ra} + \frac{2\sigma^2}{rab} + \frac{2\sigma^2}{rab} - \frac{2\sigma^2}{rab} - \frac{2\sigma^2}{rab}\right]$$

$$\mathrm{E}[SS_{A\times B}] = r\sum_{j,k}\left[(pq)_{jk} + \frac{\sigma^2}{r} - \frac{\sigma^2}{ra} - \frac{\sigma^2}{rb} + \frac{\sigma^2}{rab}\right]$$

$$\mathrm{E}[SS_{A\times B}] = ab\sigma^2 - b\sigma^2 - a\sigma^2 + \sigma^2 r\sum_{j,k}(pq)_{jk}$$

$$\mathrm{E}[SS_{A\times B}] = (a-1)(b-1)b\sigma^2 + r\sum_{j,k}(pq)_{jk}$$

**Problem 2**

a. To achieve equal variance among the estimated treatment differences, I would have one treatment of each kind in the blocks of four (O,A,B,C), one of each A,B,C in the both blocks of three, and O with each of the other treatments respectively for the final three blocks. In other words, the set of blocks would be:

$$\{\{O, A, B, C\}, \{O, A, B, C\}, \{A, B, C\}, \{A, B, C\}, \{O, A\}, \{O, B\}, \{O, C\}\}$$

This would result in the covariance matrix:

|  | (Intercept) | trtB | trtC | trtO | blk2.2 | blk2.3 | blk3.1 | blk3.2 | blk4.1 | blk4.2 |
|---|---|---|---|---|---|---|---|---|---|---|
| (Intercept) | 0.620 | -0.111 | -0.111 | -0.241 | -0.444 | -0.444 | -0.546 | -0.546 | -0.505 | -0.505 |
| trtB | -0.111 | 0.444 | 0.222 | 0.222 | -0.222 | -0.111 | -0.111 | -0.111 | -0.111 | -0.111 |
| trtC | -0.111 | 0.222 | 0.444 | 0.222 | -0.111 | -0.222 | -0.111 | -0.111 | -0.111 | -0.111 |
| trtO | -0.241 | 0.222 | 0.222 | 0.481 | -0.111 | -0.111 | 0.093 | 0.093 | 0.009 | 0.009 |
| blk2.2 | -0.444 | -0.222 | -0.111 | -0.111 | 1.111 | 0.556 | 0.556 | 0.556 | 0.556 | 0.556 |
| blk2.3 | -0.444 | -0.111 | -0.222 | -0.111 | 0.556 | 1.111 | 0.556 | 0.556 | 0.556 | 0.556 |
| blk3.1 | -0.546 | -0.111 | -0.111 | 0.093 | 0.556 | 0.556 | 0.954 | 0.620 | 0.579 | 0.579 |
| blk3.2 | -0.546 | -0.111 | -0.111 | 0.093 | 0.556 | 0.556 | 0.620 | 0.954 | 0.579 | 0.579 |
| blk4.1 | -0.505 | -0.111 | -0.111 | 0.009 | 0.556 | 0.556 | 0.579 | 0.579 | 0.808 | 0.558 |
| blk4.2 | -0.505 | -0.111 | -0.111 | 0.009 | 0.556 | 0.556 | 0.579 | 0.579 | 0.558 | 0.808 |

And give the following unscaled variances among the estimated differences:

|  | A-B | A-C | A-O | B-C | B-O | C-O |
|---|---|---|---|---|---|---|
| Unscaled Variance | 0.444 | 0.444 | 0.481 | 0.444 | 0.481 | 0.481 |

b. To get the variance estimates for the differences between 0 and the others to be $\frac{2}{3}$ that of the other differences, I used the following block assignment:

$$\{\{O, A, B, C\}, \{O, O, A, B\}, \{O, A, C\}, \{O, B, C\}, \{O, A\}, \{O, B\}, \{O, C\}\}$$

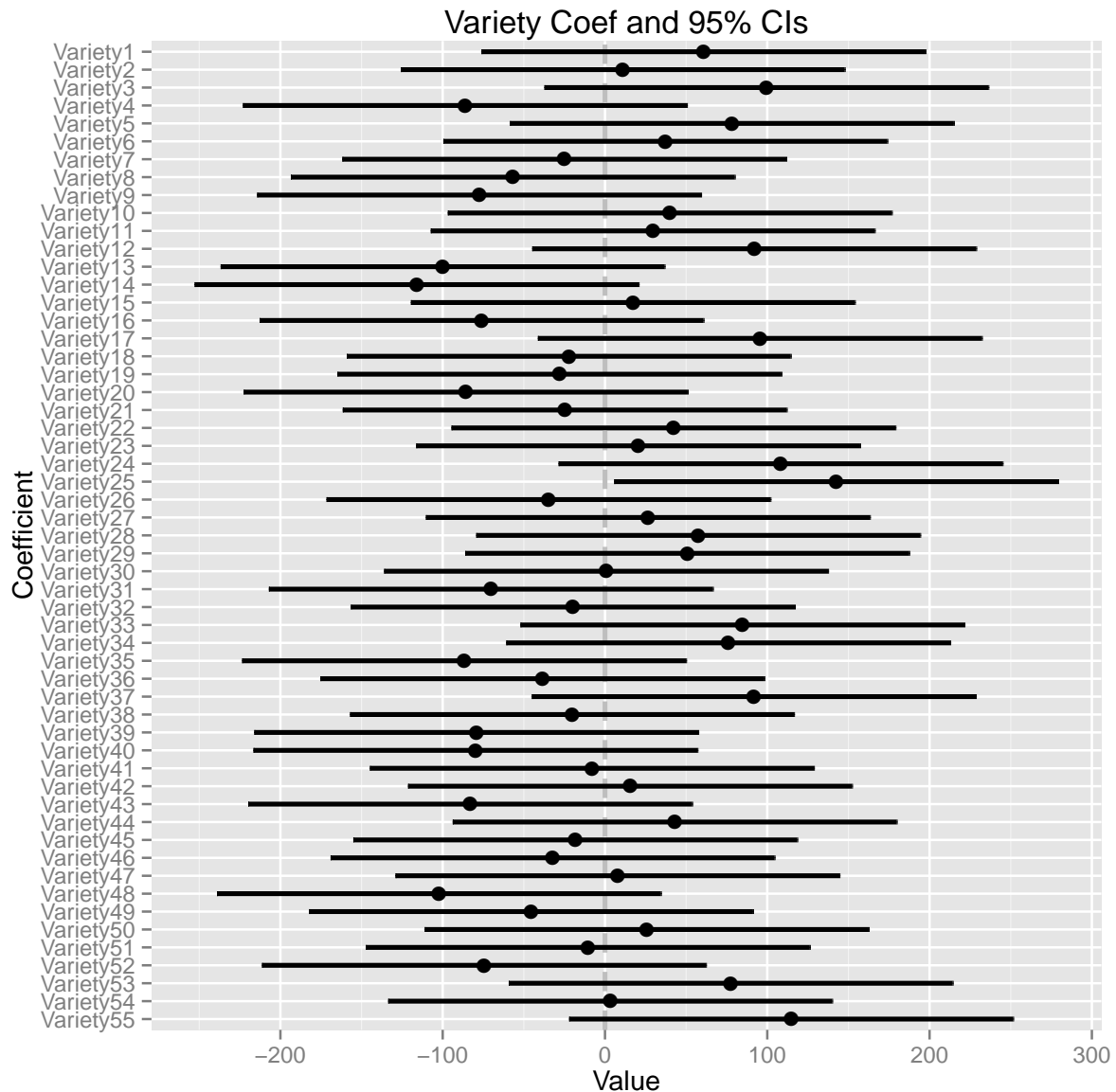This would result in the covariance matrix:

|  | (Intercept) | trtB | trtC | trtO | blk2.2 | blk2.3 | blk3.1 | blk3.2 | blk4.1 | blk4.2 |
|---|---|---|---|---|---|---|---|---|---|---|
| (Intercept) | 0.605 | -0.158 | -0.151 | -0.210 | -0.421 | -0.425 | -0.485 | -0.432 | -0.475 | -0.461 |
| trtB | -0.158 | 0.632 | 0.316 | 0.316 | -0.316 | -0.158 | -0.053 | -0.263 | -0.158 | -0.158 |
| trtC | -0.151 | 0.316 | 0.624 | 0.302 | -0.158 | -0.312 | -0.158 | -0.263 | -0.159 | -0.079 |
| trtO | -0.210 | 0.316 | 0.302 | 0.421 | -0.158 | -0.151 | -0.031 | -0.136 | -0.049 | -0.079 |
| blk2.2 | -0.421 | -0.316 | -0.158 | -0.158 | 1.158 | 0.579 | 0.526 | 0.632 | 0.579 | 0.579 |
| blk2.3 | -0.425 | -0.158 | -0.312 | -0.151 | 0.579 | 1.156 | 0.579 | 0.631 | 0.580 | 0.539 |
| blk3.1 | -0.485 | -0.053 | -0.158 | -0.031 | 0.526 | 0.579 | 0.881 | 0.565 | 0.545 | 0.513 |
| blk3.2 | -0.432 | -0.263 | -0.263 | -0.136 | 0.632 | 0.631 | 0.565 | 0.986 | 0.598 | 0.566 |
| blk4.1 | -0.475 | -0.158 | -0.159 | -0.049 | 0.579 | 0.580 | 0.545 | 0.598 | 0.817 | 0.539 |
| blk4.2 | -0.461 | -0.158 | -0.079 | -0.079 | 0.579 | 0.539 | 0.513 | 0.566 | 0.539 | 0.789 |

And give the following unscaled variances among the estimated differences:

|   | A-B | A-C | A-O | B-C | B-O | C-O |
|---|-----|-----|-----|-----|-----|-----|
| 1 | 0.632 | 0.624 | 0.421 | 0.624 | 0.421 | 0.441 |

**Problem 3**

   a. Using the standard model, variety doesn't seem to have much a significant impact on yield. I ran a linear regression with the variety and block as predictors (I used a sum constraint on variety and a treatment constraint on block). I have plotted the coefficients for variety and their confidence intervals. Since I used the sum constraint, 0 represents the average of all varieties.



As you can see, most of the confidence intervals overlap 0, indicating they are not statistically different from the mean across all varieties. One barely doesn't overlap, but the confidence intervals don't account for multiple testing, so a small number of "significant" results are to be expected. Thus, this appears to be a null plot, where variety doesn't actually predict yield significantly. To test this, I generated an ANOVA table:

|           | Df      | Sum Sq      | Mean Sq     | F value  | Pr(>F) |
|-----------|---------|-------------|-------------|----------|--------|
| Variety   | 55.000  | 954994.888  | 17363.543   | 0.875    | 0.712  |
| Block     | 3.000   | 723630.442  | 241210.147  | 12.162   | 0.000  |
| Residuals | 165.000 | 3272436.308 | 19832.947   |          |        |

Indeed, the F test on Variety yields a high p-value of .712. We can not reject the null hypothesis that different varieties have the same expected yield.
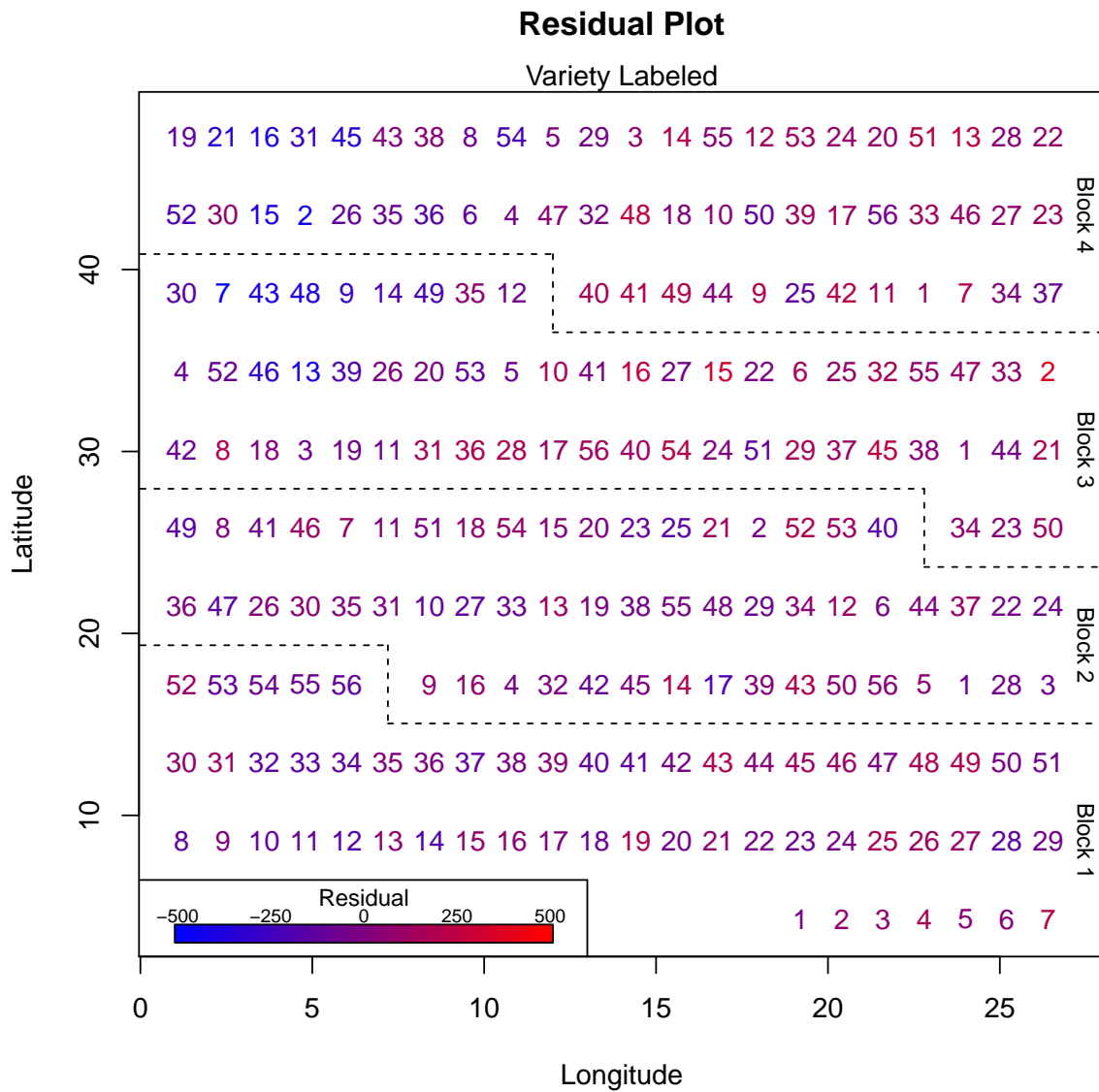
b. I calculated this estimate as a contrast on the estimated parameters. Where $\beta_i$ was the estimated effect for variety $i$,

$$a_i = \begin{cases} \frac{1}{20} & \text{if } i \leq 20 \\ \frac{1}{36} & \text{if } i > 20 \end{cases}$$

Thus, my estimate was $d = \sum_{i=1}^{56} a_i \beta_i = -8.813$. Then, my confidence interval was gotten in the usual way for a linear combination of parameters from a linear regression:

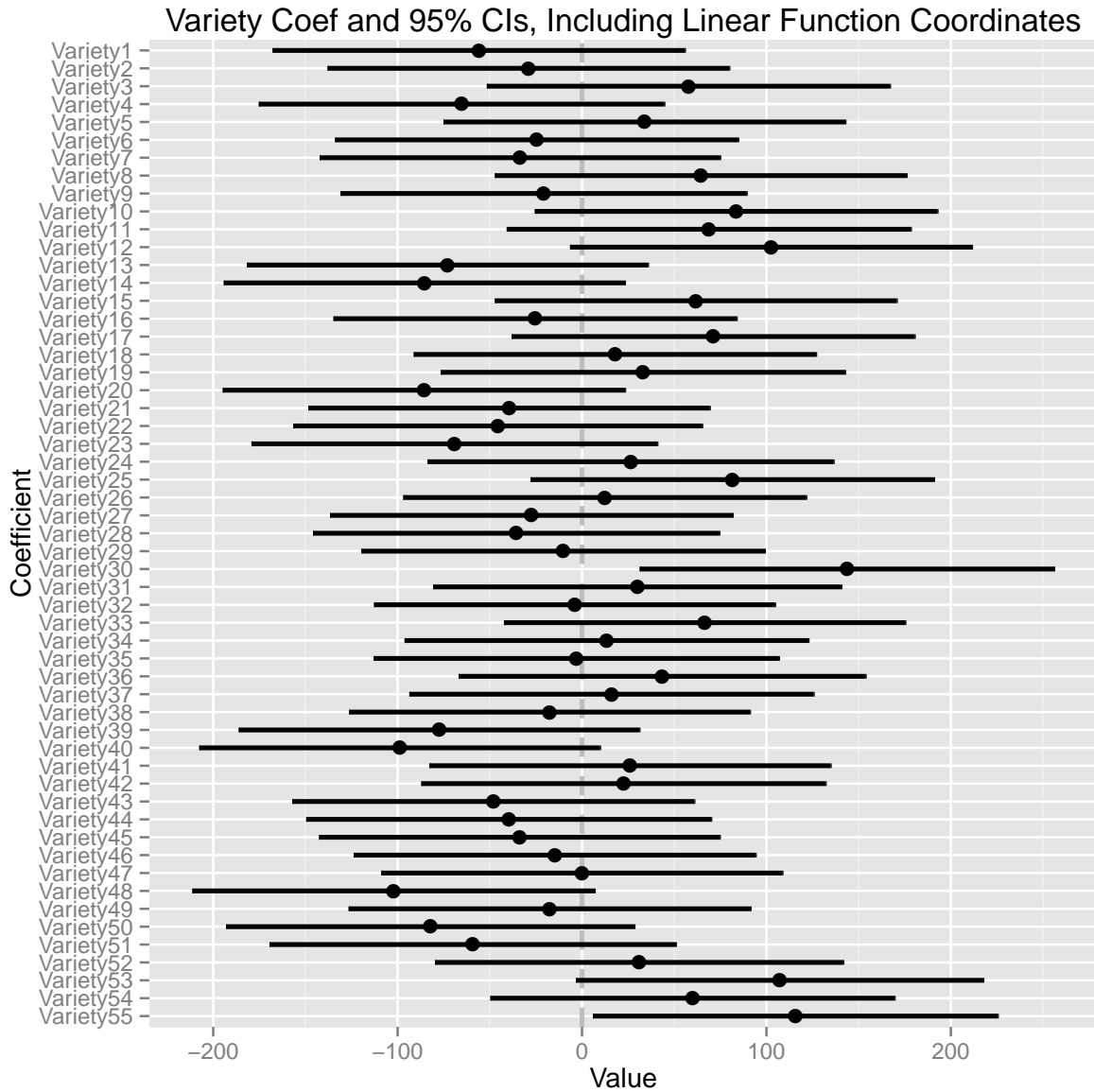$$d \pm t_{165}^{.025} \sqrt{a^t \Sigma a} = [-47.586, 29.961]$$

c. I have plotted the residuals below:
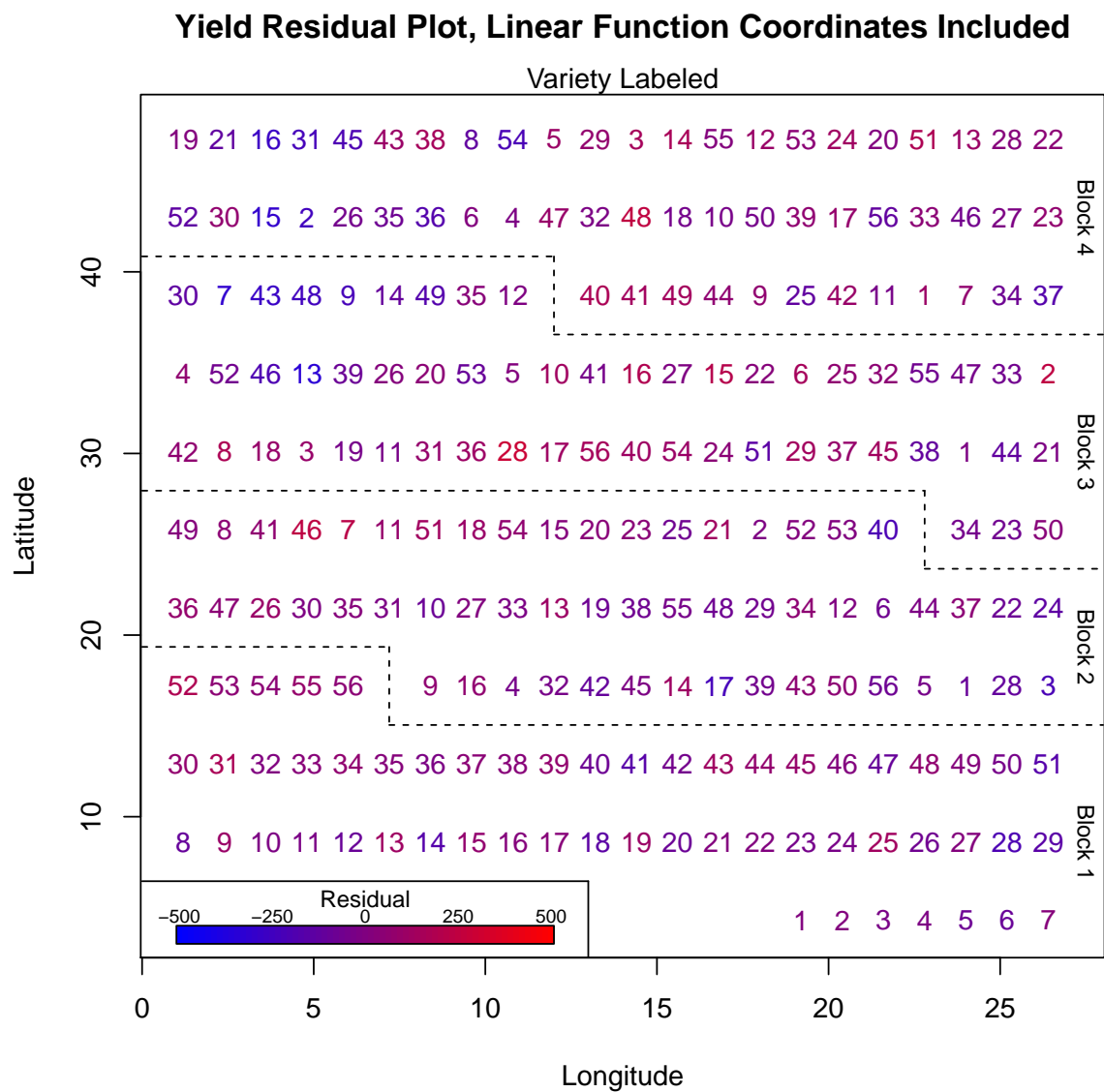


**Residual Plot**

It appears that the units in the top left corner (high latitude, low longitude), tend not to have mean 0, instead having negative mean. This indicates that the blocking didn't really account for the variation in yield due to geography as the study assumed. Any of the varieties with one plot in that region would be adversely affected, but those that have a plot in both block 3 and 4 which fall in that region would be the

6

worst affected, driving down their predicted yield. This seems to include varieties 30, 26, 52, 4, and possibly a few others.

d. After including longitude and latitude as predictors in the model, we see some movement in the predicted effects of the varieties:
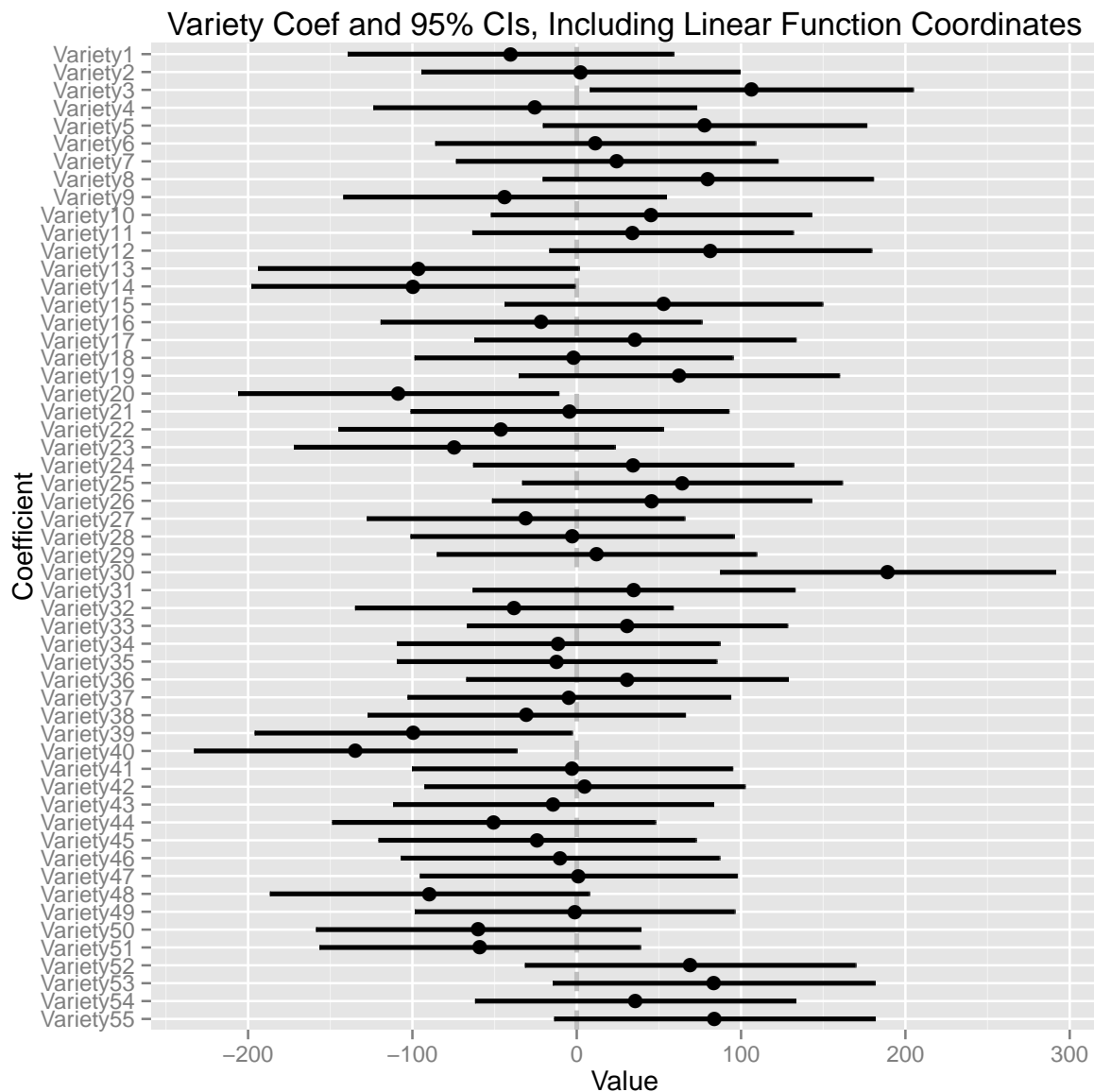


From the previous regression, most of the predicted effects stay approximately the same. Some, such as 1, flip signs, though their CI still cover 0. A few, such as 55 and 53 are nudged over enough so that their CI don't cover 0, but with so many tests that distinction isn't meaningful. The only meaningful change I see is that 30 made a large shift, becoming the most significant effect while before the estimate was very close to 0. Looking at the residuals again, we still see somewhat of the same problem we had before, though less substantial:
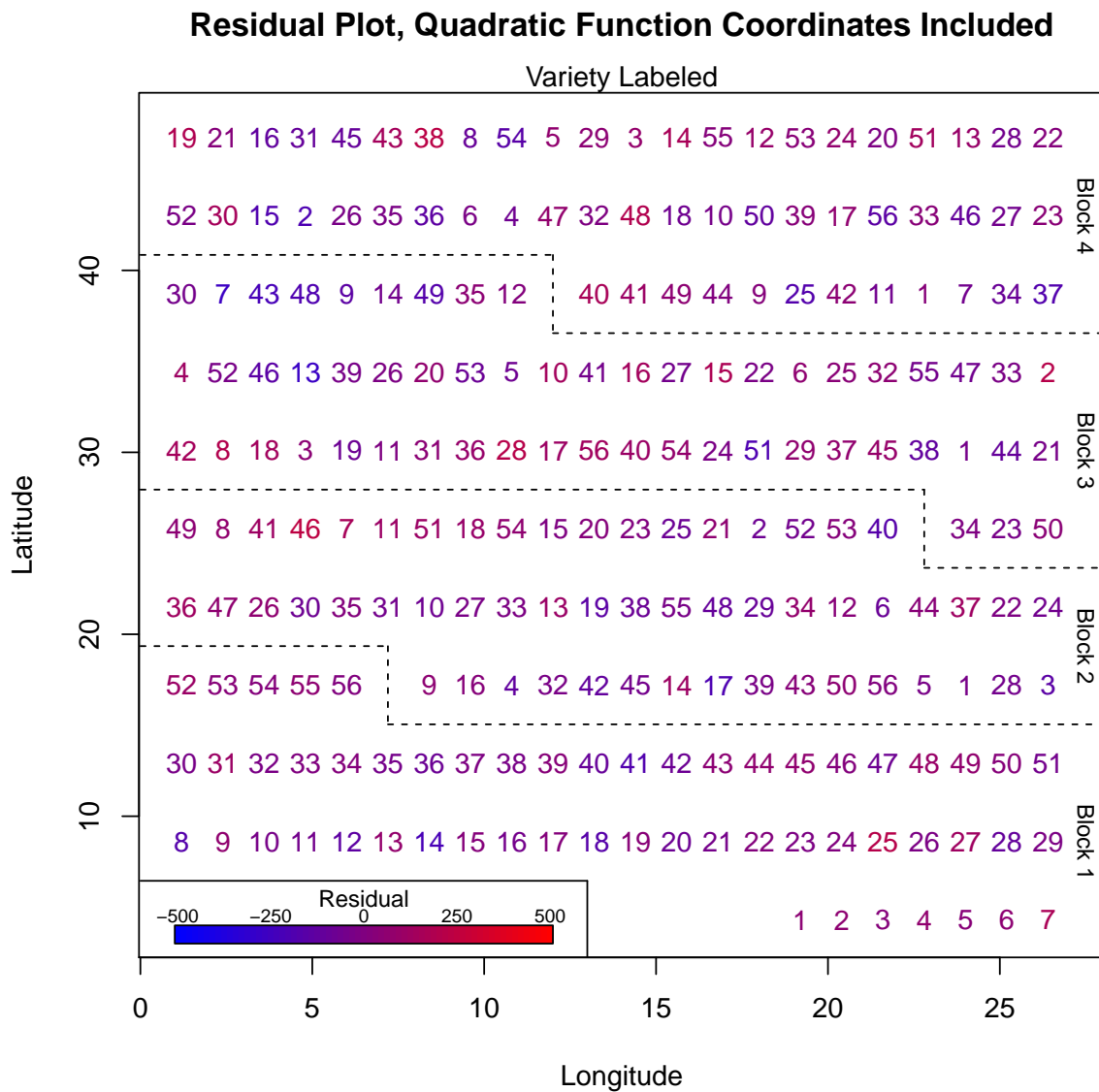
**Yield Residual Plot, Linear Function Coordinates Included**

Variety Labeled



There still remains what looks like a meaningful dip in yield in that top left corner.

e. With the quadratic function of the coordinates, we finally see substantial differentiation between the varieties:

Variety Coef and 95% CIs, Including Linear Function Coordinates

Variety 30 appears to be quite distinct from the rest of the varieties; it's p-value is .000368, which is sufficiently small that we can reject the null assumption at an alpha of 5% even with a Bonferoni adjustment. A few other varieties' CIs now no longer cover 0, such as 39, 40, and 20, but they aren't by very much. Since I am using the sum constraint, this may be a function of variety 30 dragging up the average rather than than these varieties being meaningfully worse than all the varieties excluding 30. Looking at the residuals again, we seem to have finally solved the issue with the top left corner:

## Residual Plot, Quadratic Function Coordinates Included



By eyeball, I can not say that that plots in that region do worse than the rest of the plots, and no other area jumps out as doing particularly well or poorly. If there is a lingering geographic effect, it does not seem to be substantial.

f. I do not think this study design was well chosen. The intent was clearly that the blocks would account for heterogeneity in yield over different areas of the farm, but this did not hold true in practice. Additionally, the blocks allowed replicates of the same variety to be planted relatively closely in the field. For instance, two of variety 30's replicates were right next to each other in the top right corner, and all of 30's replicates were in the four left most columns of plots. A better solution would have been to create smaller blocks which had much closer to homogeneous growing conditions, and then assign varieties to blocks such that they are more evenlyh spread over the field. For instance, the field might have been broken up into approximately a $4 \times 4$ grid, where each replicate of a particular variety would be placed in a block such that none of the other replicates of that variety would be in the same row or column of blocks. This would have prevented them being overly clustered together. The cost of this would be that the experiment would no longer have a full factorial design, but with 14 treatments in each block there would still be substantial grounds to compare varieties.