# 1 growth

We fit a second-order model using the following equation, check the summary and analysis of variance.

```
> h<-lm(Yield~x1+x2+x3+I(x1^2)+I(x2^2)+I(x3^2)+x1*x2+x2*x3+x1*x3,data = growth)
> summary(h)
> pure.error.anova(h)
```

The summay and anova table are shown as shown in figure 1 and figure 2

```
> summary(h)

Call:
lm(formula = Yield ~ x1 + x2 + x3 + I(x1^2) + I(x2^2) + I(x3^2)
 +
    x1 * x2 + x2 * x3 + x1 * x3, data = growth)

Residuals:
    Min      1Q   Median      3Q     Max
-15.6661  -9.1577  -0.6661   9.1718  17.3339

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  100.666      5.564  18.093  5.7e-09 ***
x1             1.271      3.691   0.344  0.73765
x2             1.361      3.691   0.369  0.71998
x3            -1.494      3.691  -0.405  0.69411
I(x1^2)       -3.767      3.593  -1.048  0.31912
I(x2^2)      -12.430      3.593  -3.459  0.00613 **
I(x3^2)       -9.601      3.593  -2.672  0.02342 *
x1:x2          2.875      4.823   0.596  0.56436
x2:x3         -4.625      4.823  -0.959  0.36020
x1:x3         -2.625      4.823  -0.544  0.59819
---
Signif. codes:
0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 13.64 on 10 degrees of freedom
Multiple R-squared:  0.6631,    Adjusted R-squared:  0.3598
F-statistic: 2.186 on 9 and 10 DF,  p-value: 0.1194
```

Figure 1: summary of second order model

```
> pure.error.anova(h)
Analysis of Variance Table

Response: Yield
           Df  Sum Sq Mean Sq F value  Pr(>F)
x1          1   22.08   22.08  0.1285 0.73467
x2          1   25.31   25.31  0.1473 0.71695
x3          1   30.50   30.50  0.1775 0.69105
I(x1^2)     1   46.37   46.37  0.2698 0.62564
I(x2^2)     1 1916.92 1916.92 11.1535 0.02056 *
I(x3^2)     1 1328.46 1328.46  7.7296 0.03889 *
x1:x2       1   66.12   66.12  0.3847 0.56225
x2:x3       1  171.13  171.13  0.9957 0.36417
x1:x3       1   55.13   55.13  0.3207 0.59564
Residuals  10 1860.95  186.09
 Lack of fit 5 1001.61  200.32  1.1656 0.43529
 Pure Error  5  859.33  171.87
---
```

Figure 2: anova of second order model

The p-value of lack of fit is 0.43529, which is greater than 0.05. Thus, we could conclude that the second model is adequate to represent the data. Fitted model: $\hat{y} = 100.666 + 1.271\hat{x}_1 + 1.361\hat{x}_2 - 1.494\hat{x}_3 - 3.767\hat{x}_1{}^2 - 12.430\hat{x}_2{}^2 - 9.601\hat{x}_3{}^2 + 2.875\hat{x}_1\hat{x}_2 - 4.625\hat{x}_2\hat{x}_3 - 2.625\hat{x}_1\hat{x}_3$

$$B_{3,3} = \begin{pmatrix} -3.7670 & 1.4375 & -1.3125 \\ 1.4375 & -12.4300 & -2.3125 \\ -1.3125 & -2.3125 & -9.6010 \end{pmatrix}, b = \begin{pmatrix} 1.271 \\ 1.361 \\ -1.494 \end{pmatrix}, x_s = -\tfrac{1}{2}B^{-1}b =$$

$$\begin{pmatrix} 0.260 \\ 0.111 \\ -0.140 \end{pmatrix},$$ and the eigenvalues of matrix B are $$\begin{pmatrix} -3.078495 \\ -8.953226 \\ -13.766279 \end{pmatrix}.$$ All eigen-values are negative, which makes sure that $X_s$ is the maximum point.

## 2 average age

**a)** The sample design is simple random sampling without replacement. Under SRSWOR, the sample mean $\bar{y}$ is an unbiased estimator of $\bar{Y}$, thus the estimator of mean age for children is $\bar{y} = \frac{9*13+10*35+11*44+12*69+13*36+14*24+15*7+16*3+17*2+18*5}{240} = 12.08$ The $v(\bar{y})$ is an unbiased estimator of $V(\bar{y})$, and $v(\bar{y}) = \frac{s^2}{n} = \frac{3.705}{240} = 0.015$, thus, the standard error $se(\bar{y}) = \sqrt{v(\bar{y})} = 0.124$. And the 95% confidence interval for the average age is $\bar{y} \pm Z_{\alpha/2}s\sqrt{\frac{1}{n}} = 12.08 \pm 0.243$ **b)** We determine the sample size based on this formula: $n = \frac{Z_{\alpha/2}{}^2 S^2}{e^2} = \frac{1.96^2 * 3.705}{0.5^2} = 56.93$, hence, the minimun sample size is 57.

## 3 clams

First, we calculate $N_h h = 1,,,4$, $N_1 = 222.81 * 25.6 = 5704$, $N_2 = 49.61 * 25.6 = 1270$, $N_3 = 50.25 * 25.6 = 1287$, $N_4 = 197.81 * 25.6 = 5064$, $N = N_1 + N_2 + N_3 + N_4 = 13325$. Then we obtain the $\bar{y}_{st} = \sum_{h=1}^{H} W_h \bar{y}_h = 1.36$. After that, we can have the estimator of the total number of bushels $\hat{t}_{st} = N\bar{y}_{st} = 13325 * 1.36 = 18122$.

The variance of $\hat{y}_{st}$: $v(\hat{y}_{st}) = \sum_{h=1}^{H} W_h{}^2(1 - n_h/N_h)s_h{}^2/n_h = 0.0327$, thus, the variance of $\hat{t}_{st}$: $v(\hat{t}_{st}) = N^2 v(y_{st}) = 13325^2 * 0.0327 = 5806069$ the standard error is $se(\hat{t}_{st}) = \sqrt{v(\hat{t}_{st})} = 2410$

## 4 totoal number of acres

**a)** Use ratio estimation to estimate the total number of acres:

$\hat{R} = \frac{\bar{y}}{\bar{x}} = \frac{mean(acres92)}{mean(farms87)} = 459.8975$

$\hat{t_{yr}} = \hat{R}t_x = 459.8975 * 2087759 = 960,155,061$

**b)** Use the regression estimation:

$\hat{\beta}_0 = 263098.45, \hat{\beta}_1 = 58.09$

$\hat{\bar{y}}_{req} = \hat{\beta}_0 + \hat{\beta}_1\bar{x} = 263098.45 + 58.09 * 2087759/3078 = 302,500$

$\hat{t}_{yreq} = N\hat{\bar{y}}_{req} = 3078 * 302500 = 931,095,000$

**c)** In order to find the method with most precision, we calculate the standard variance of $\hat{t}_y$.

ratio estimation with auxiliary variable acres87, $se(\hat{t}_{yra87}) = \sqrt{var(\hat{t}_y)} =$

$\sqrt{N^2(1 - \frac{n}{N})\frac{1}{n}\frac{1}{n-1}\sum_{i\in s}(y_i - \hat{R}x_i)^2} = 5,344,567$

ratio estimation with auxiliary variable farms87,

$se(\hat{t}_{yrf87}) = \sqrt{var(\hat{t}_y)} = \sqrt{N^2(1 - \frac{n}{N})\frac{1}{n}\frac{1}{n-1}\sum_{i\in s}(y_i - \hat{R}x_i)^2} = 65,364,822$

regression estimation with auxiliary variable farms87,

$se(\hat{t}_{yregf87}) = \sqrt{var(\hat{t}_y)} = \sqrt{N^2(1 - \frac{n}{N})\frac{1}{n}\frac{1}{n-1}\sum_{i\in s}(y_i - \beta_0 - \beta_1 * x_i)^2} =$
$58,065,813$

Based on the variances, we can tell that ratio estimation has the most precision since its variance is minimum among these three methods.

# 5 Neyman allocation

a)

$$V_{Neyman}(\hat{t}_{str}) = N^2 V(\bar{y}_{st}) = N^2 \sum_{h=1}^{H} W_h{}^2 (1 - \frac{n_h}{N_h}) S_h{}^2 / n_h$$

$$= \sum_{h=1}^{H} N_h{}^2 (1 - \frac{n_h}{N_h}) S_h{}^2 / n_h$$

$$= \sum_{h=1}^{H} N_h{}^2 (1 - \frac{\frac{N_h S_h n}{\sum_{h=1}^{H} N_l S_l}}{N_h}) S_h{}^2 \frac{\sum_{h=1}^{H} N_l S_l}{N_h S_h n}$$

$$= \sum_{h=1}^{H} N_h S_h (1 - \frac{S_h n}{\sum_{h=1}^{H} N_l S_l}) \frac{\sum_{h=1}^{H} N_l S_l}{n}$$

$$= \sum_{h=1}^{H} N_h S_h (\frac{\sum_{h=1}^{H} N_l S_l}{n} - S_h)$$

$$= \frac{1}{n} \sum_{h=1}^{H} N_l S_l \sum_{h=1}^{H} N_h S_h - \sum_{h=1}^{H} N_h S_h{}^2$$

$$= \frac{1}{n} (\sum_{h=1}^{H} N_h S_h)^2 - \sum_{h=1}^{H} N_h S_h{}^2$$

b)

$$V_{prop}(\hat{t}_{str}) - V_{Neyman}(\hat{t}_{str})$$

$$= \frac{N}{n} \sum_{h=1}^{H} N_h S_h{}^2 - \sum_{h=1}^{H} N_h S_h{}^2 - \frac{1}{n} (\sum_{h=1}^{H} N_h S_h)^2 + \sum_{h=1}^{H} N_h S_h{}^2$$

$$= \frac{N}{n} \sum_{h=1}^{H} N_h S_h{}^2 - \frac{1}{n} (\sum_{h=1}^{H} N_h S_h)^2$$

$$= \frac{N^2}{n} \sum_{h=1}^{H} \frac{N_h}{N} S_h{}^2 - \frac{N^2}{n} (\sum_{h=1}^{H} \frac{N_h}{N} S_h)^2$$

$$= \frac{N^2}{n} [\sum_{h=1}^{H} \frac{N_h}{N} S_h{}^2 - (\sum_{h=1}^{H} \frac{N_h}{N} S_h)^2]$$

$$= \frac{N^2}{n} \sum_{h=1}^{H} \frac{N_h}{N} (S_h - \sum_{l=1}^{H} \frac{N_l}{N} S_l)^2$$