

PREDICTING APPLE SIZE: PRELIMINARY RESULT

BU-201-M

O. Ladipo and A. Hedayat

July, 1965

Abstract

The final size of an apple, that is to say the size of an apple at harvest, does not depend only on its size four weeks before harvest, but it also depends on a number of factors which, acting simultaneously, determine how big or small the apple is. In order to predict apple production more accurately several factors must be considered.

---

Biometrics Unit, Cornell University, Ithaca, New York

## PREDICTING APPLE SIZE: PRELIMINARY RESULT

BU-201-M

O. Ladipo and A. Hedayat

July, 1965

### Nature of the Preliminary Problem:

The question here is "Can we predict final apple size on the basis of one factor, namely the size of the apple in mid-August?". Our first reaction is that final apple size is a response which is affected by a number of quantitative factors. Which of these factors are the most efficient in predicting final apple size? From here on we will refer to final apple size as  $Y$ ; the other factors will be as follows:

$X_1$  = Apple size as of (or near) August 10

$X_2$  = Apple size as of (or near) August 21

$X_3$  = Nitrogen level as of August 15

$X_4$  = Total number of dry days

$X_5$  = Number of dry days before August 15

$X_6$  = Crop-load

The ultimate purpose of this study is to predict total production of apples in New York State as early as possible before harvesting time. At this stage we want to show that several factors affect production. After finding those factors that are most efficient in predicting apple size we will then aim at constructing a formula (a model) and designing a sampling procedure both of which, we hope, can be used annually before harvesting to predict what the year's total apple production will be.

The Data:

The data we are working with in this preliminary stage were collected at the Geneva Experimental Station. The set of data consists of observations in six orchards. Each observation is "fruit volume in cubic inches". The measurements were taken at intervals of approximately ten days starting July 2<sup>nd</sup> and over a period of four years -- 1961-1964.

It is our feeling that the data were collected as a matter of routine. By this we mean they lacked any design as to the purpose of the collection. We do not know whether the value 6.13 (for August 21<sup>st</sup>, orchard No. 1, for 1962) is the average volume of all apples for that orchard or a sample mean volume. We, for the purpose of the present problem, regard it as a mean volume of the apples in a particular orchard. Under "Soil Moisture" it is difficult to decide what "Dry days" really mean. Here we assume, following the footnotes on the data sheet, that the figures represent numbers of days with at least minimum amount of moisture. We scored crop-load as light = 1, moderate = 2, and heavy = 3.

We do feel that temperature is an important variable which should have been included. Maybe the pomologists can help us with an explanation for its omission. Other things we would have wanted to see are rainfall (in addition to dry days), thinning of apples and age of orchard. With sufficient information we can find the yield per acre and knowing the number of acres per orchard we will be able to establish the production of apples in an orchard.

Analysis:

There is no gainsaying the fact that final apple size (i.e., its size in mid-September) is dependent on its size in mid-August. This is justified by

the high correlation coefficient  $r_{yx_1}$ , which is .88. The question, however, is: Can we do better in the prediction of Y if we consider other factors?.

In answer to this we resorted to multiple linear regression (see diagrams I-V, to each of which a regression line can be fitted). Setting up the equation we have:

$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6$$

We claim that adding more variables may give us more information. To check this we present Table 1.

Table 1

Variables added in order	Meaning of variable	R	R <sup>2</sup>	Additional information gained	Remarks
X <sub>1</sub>	size as of 8/15	.88	.776	-	Almost 23% of variability in Y is not accounted for by regression
X <sub>4</sub>	Total dry days	.9129	.8334	.0574	Additional information gained by considering these three factors is 13.85% and only 8.6% of variability is not accounted for by regression
X <sub>3</sub>	Nitrogen level	.9440	.8912	.0578	
X <sub>6</sub>	Crop-load	.9563	.9145	.0233	
X <sub>2</sub>	Size as of 8/21	.9614	.9242	.0097	Additional gain is only 1.09%
X <sub>5</sub>	Dry days before 8/15	.9620	.9254	.0012	

From the above table we observe that not all factors contribute a lot of information when added. The bulk of the extra information comes from fixing  $X_4$  and  $X_3$  which yields 11.5% more information while if  $X_2$  and  $X_5$  are fitted along with  $X_1$ ,  $X_3$ ,  $X_4$  and  $X_6$  we only gained 1.01% more information. It seems, therefore, not worth while to include  $X_2$  and  $X_5$ . Our equation then becomes

$$Y = \alpha + \beta_1 X_1 + \beta_3 X_3 + \beta_4 X_4 + \beta_6 X_6$$

If, however, we can take two measurements of fruit volume the addition of  $X_2$  will give us a little additional information and our regression equation will then be

$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_6 X_6$$

Based on the crude data at our disposal we regard this as our optimal equation. The inclusion of  $X_5$  adds only 0.12% to the amount of information. Some of the combinations with the amount of information are given in Table 2. In Table 3 we try to show that though 92.54% of variability in  $Y$  may be accounted for by regression, yet the individual coefficient of regression ( $b_i$ ) need not be significant at 5% level. In Table 3 we find that, using a t-test not all of the  $b_i$ 's are significant at 5% level.

In order to test for the relative importance of each of the independent variables we considered  $\beta'_i$  "since each  $\beta'$  is independent of the original units of measurement, a comparison of any two indicates the relative importance of the independent variables involved". (Steel and Torrie, 1960, p. 284.) In Table 3 we showed the relative importance of each factor in predicting  $Y$ .

Table 2

Variables in equation	R	R <sup>2</sup>	Variables left out
X <sub>1</sub> , X <sub>2</sub> , X <sub>4</sub> , X <sub>5</sub> , X <sub>6</sub>	.94	.88	X <sub>3</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>3</sub> , X <sub>4</sub> , X <sub>5</sub>	.9565	.9149	X <sub>6</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>3</sub> , X <sub>5</sub> , X <sub>6</sub>	.9558	.91355	X <sub>4</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>3</sub> , X <sub>4</sub> , X <sub>6</sub>	.9614	.9242	X <sub>5</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>3</sub> , X <sub>6</sub>	.9155	.84	X <sub>4</sub> , X <sub>5</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>4</sub> , X <sub>5</sub>	.9373	.8785	X <sub>3</sub> , X <sub>6</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>4</sub>	.9138	.8350	X <sub>3</sub> , X <sub>5</sub> , X <sub>6</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>5</sub>	.8967	.8041	X <sub>3</sub> , X <sub>4</sub> , X <sub>6</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>6</sub>	.8816	.7772	X <sub>3</sub> , X <sub>4</sub> , X <sub>5</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>3</sub>	.9106	.8292	X <sub>4</sub> , X <sub>5</sub> , X <sub>6</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>3</sub> , X <sub>5</sub>	.9424	.8882	X <sub>4</sub> , X <sub>6</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>3</sub> , X <sub>4</sub>	.95235	.90696	X <sub>5</sub> , X <sub>6</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>5</sub> , X <sub>6</sub>	.8983	.8069	X <sub>3</sub> , X <sub>4</sub>
X <sub>1</sub> , X <sub>2</sub> , X <sub>4</sub> , X <sub>6</sub>	.91503	.83729	X <sub>3</sub> , X <sub>5</sub>

Nitrogen level is the most consistent. The importance of the other factors (as shown by their ranks) changes with respect to the combination. If we leave out X<sub>4</sub> and replace it with X<sub>5</sub> we lose only 1.06% of the information. We will be forced to do this because, in predicting final apple size, the total number of dry days will not be available to us, and we will, therefore, use number of dry days before August 15<sup>th</sup>. It has been stated in the letter of March 11, 1965 from Professor Forshey to Professor Federer that "at least 30 days of deficient soil moisture are required for a significant reduction in fruit size". If this

Table 3

Combination of variables	$b_1$	Value of $t_1$ for testing $b_1$	Level of $t_1$ being significant	$b_1'$	Rank of $X_1$ in predicting -Y	R	$R^2$	Remarks
1 Size at 8/10	+0.384042	0.355	> 50%	+0.232962	4	0.962	0.9254	
2 Size at 8/21	+0.476965	0.553	> 50%	+0.321627	3			
3 Nitrogen	+3.50273	2.379	3%	+0.423022	2			
4 Total days	-0.020606	1.195	25%	-0.443585	1			
5 Day to 8/15	+0.010286	0.373	> 50%	+0.144365	5			
6 Crop-load	-0.20787	1.127	25%	-0.143572	6			
Constant term	-2.10978							
1	+0.889263	4.074	1%	+0.059346	5	0.961	0.9242	Deleting variable five or total dry days before 8/15 does not decrease the amount of information for predicting Y.
2	+0.693389	1.135	25%	+0.467566	1			
3	+3.383050	3.387	1%	+0.462607	2			
4	-0.014381	3.371	1%	-0.309577	3			
6	-0.238530	1.510	15%	-0.164748	4			
Constant term	-6.70803							
1	-0.499068	0.619	> 50%	-0.302737	3	0.956	0.914	Here $b_1$ has negative sign that should be interpreted, because we expect it should have always positive sign.
2	+1.162220	1.767	15%	+0.783708	1			
3	+4.450220	3.512	1%	+0.537449	2			
5	-0.021419	2.954	1.5%	-0.301488	4			
6	-0.295910	1.713	15%	-0.204421	5			
Constant term	-3.63809							
1	+0.889263	4.074	1%	+0.53931	1	0.956	0.914	Again here it shows that deleting $X_5$ does not decrease the amount of information, and by these four variables we can predict Y as close as with 6 variables in case one.
3	+3.192230	1.135	25%	+0.385524	2			
4	-0.014794	3.387	1%	-0.318462	3			
6	-0.272238	3.371	1%	-0.188029	4			
Constant term	-1.07596	1.510	15%					

statement is true we would expect a significant difference between the means of two groups -- apple size in orchards with more than 30 days of deficient soil moisture and apple size in orchards with less than 30 days of deficient soil moisture. We tested this and found that the difference is not significant. The test is as follows:

$$\bar{X}_1 = 7.96 \text{ mean for orchards with } > 30 \text{ dry days}$$

$$\bar{X}_2 = 8.40 \text{ mean for orchards with } < 30 \text{ dry days}$$

$$s_p^2 = (SS_{X_1} + SS_{X_2}) / (n_1 + n_2 - 2) \\ = 19/14$$

$$s_d = \sqrt{\frac{19}{14} \left( \frac{16}{63} \right)} = .624$$

$$t = .72 / .624 = 1.153 \text{ with } 14 \text{ d.f.}$$

$$t < t_{.05, 14} = 2.145$$

Conclusion:

Difference,  $\bar{X}_2 - \bar{X}_1$ , is not significant. If we, however, consider this factor (dry days before August 15<sup>th</sup>) with some other factor, say Nitrogen level, we might reach a different conclusion if our earlier conclusion is wrong and the latter right. To check this we therefore examined the regression of final apple size in each of the two groups on Nitrogen to see if there is any significant difference in the "groups" adjusted means -- 7.5734 and 8.5431. With an  $F = 3.0877 < F_{.05}(1, 12) = 4.75$  we conclude that our earlier conclusion holds and that the statement in the letter needs some explanation.



Summary of calculations for the above follows.

$$\hat{U}_{y_1 \cdot \bar{X}} = \bar{Y}_1 - b_1(\bar{X} - \bar{\bar{X}}) = 7.5734$$

$$\hat{U}_{y_2 \cdot \bar{X}} = \bar{Y}_2 - b_2(\bar{X} - \bar{\bar{X}}) = 8.4531$$

$$H_0 : \hat{U}_{y_1 \cdot \bar{X}} - \hat{U}_{y_2 \cdot \bar{X}} \neq 0$$

Pooled residual SS = 10.23 with 12 d.f.

Sum of squares for  $\bar{Y}_{\bar{X}} = 2.6246$  with 1 d.f.

$$F = \frac{2.6246/1}{10.23/12} = \frac{12(2.6246)}{10.23} = 3.0877$$

$$F_{.05}(1,12) = 4.75$$

Conclusion:

$H_0$  is rejected.

Anticipating the question of why we fitted a linear rather than a cur-linear regression we constructed diagrams I-V. We can fit a linear regression to each one of the five and this justifies the fitting of multiple linear regression.

In all the above analyses we ignored the year effects. The available data for each year are small and with small samples we have to be cautious about any conclusions, tests of significance and be a little skeptical, too (Snedecor, 1956). Our reluctance to use small samples (each set of annual data) is based on the fact that with five independent variables ( $X_1, X_2, X_3, X_5, X_6$ ) and only six observations we have zero degrees of freedom for error ( $n-k-1 = 6-5-1$ ),

whereas when we pooled all annual data (1961 - 1964) we have 16 observations and hence 10 degrees of freedom for error in the multiple regression analysis of variance. Table 4 shows the different F's for each variable after the others have been fitted. This is done by the use of step-wise multiple regression analysis (Searle and Primer, 1964).

Table 4

Variable being fitted	F, after preceeding variables are fitted	Error degrees of freedom	Standard error of Y	Remarks: (testing at $\alpha = .05$ )
$X_1$	48.61	14	0.58	Significant
$X_4$	4.45	13	0.52	Not significant
$X_3$	6.36	12	0.44	Significant
$X_6$	3.00	11	0.404	Not significant $F(1,11)$
$X_2$	1.29	10	0.40	Not significant
$X_5$	0.14	9	0.42	Not significant

A person might expect that the standard error of Y has something questionable, especially the last three values. Recall from Table 1 that after fitting  $X_1$ ,  $X_3$ ,  $X_4$ ,  $X_6$  and  $X_2$ , only 0.12% of variability in Y is explained by fitting  $X_5$ . This small amount of reduction in the variability of Y does not seem to be worth losing one degree of freedom, and this explains why the last standard error of Y in Table 4 is greater than the preceding one. The degrees of freedom degrees by 10% (i.e. from 10 to 9), while the variability degrees only by 0.12%.

It is easy to see from the above that the most important factors are size as of August 15<sup>th</sup> (10-15), nitrogen level, while  $X_4$  - total dry days is close

to being significant because  $F_{.05}(1,13) = 4.67$ . This is the same type of conclusions we reached via Table 3.

The analysis of variance for all the combinations mentioned so far are available in the computer sheets.

As a result of the scatter diagrams we feel we might do better under a log transformation. Suppose we set up a regression formula of the form

$$C = \alpha^* X_1^{\beta_1} X_2^{\beta_2} X_3^{\beta_3} \dots X_k^{\beta_k}$$

where  $\alpha^*$ ,  $\beta_1$ ,  $\beta_2$ , ...,  $\beta_k$  are the constants to be estimated and  $X_1, X_2, \dots, X_k$  are observed variables, then the transformation is:  $Y = \log C$  if  $\log \alpha^* = \alpha$  then

$$Y = \alpha + \beta_1 \log X_1 + \beta_2 \log X_2 + \dots + \beta_k \log X_k .$$

We hope to report on this as the analyses come out of the computer.

#### LITERATURE CITED

1. Steel, R. G. D., and Torrie, J. H. Principles and Procedures of Statistics. McGraw-Hill, New York, 1960.
2. Snedecor, G. W. Statistical Methods. (Fifth edition). Iowa State Univ. Press, Ames, Iowa, 1956.
3. Searle, S. R., and Primer, Mrs. P. L. Multiple Regression Analysis. (Third Issue), 1964.

Diagram I: Final Size  $\times$  Total Moisture (Dry Days)

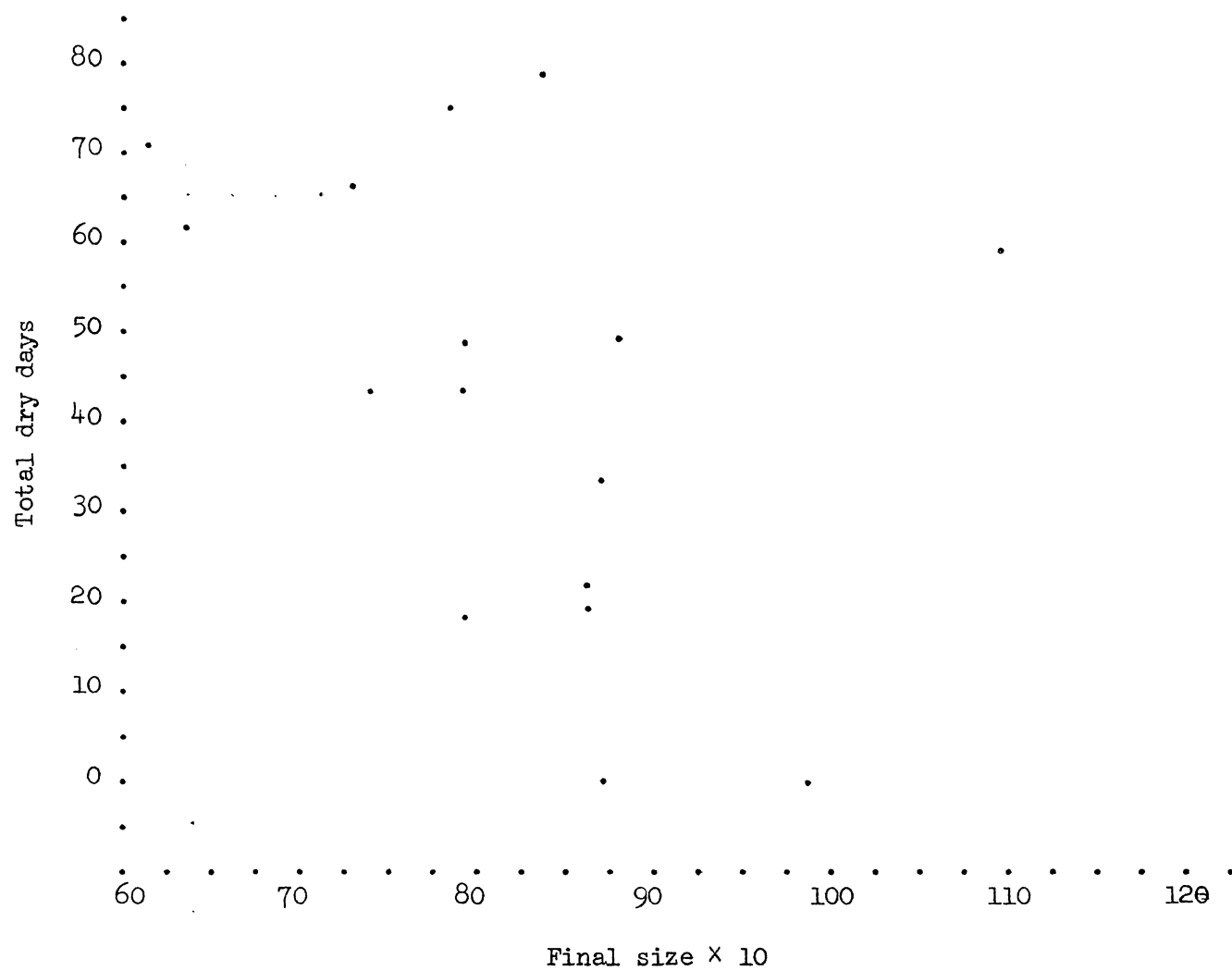
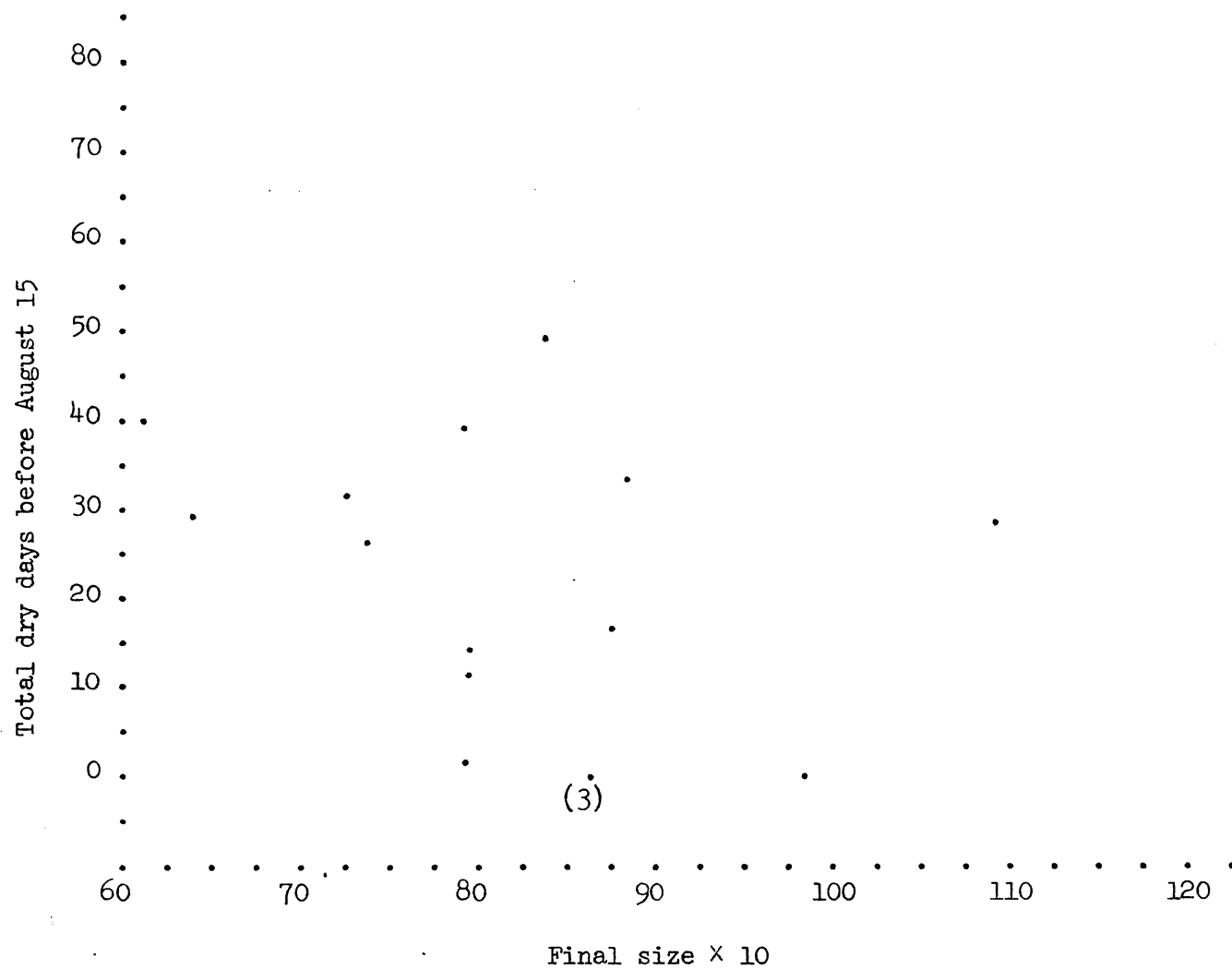


Diagram II: Final Size  $\times$  Number of Dry Days Before August 15



Scatter plot showing the relationship between Final apple Size (X-axis) and Final apple Weight (Y-axis). The X-axis ranges from 6 to 11, and the Y-axis ranges from 15 to 21. The data points show a positive correlation, indicating that as the final apple size increases, the final apple weight also tends to increase.

Final apple Size	Final apple Weight
6.0	15.5
6.1	15.8
6.3	18.0
6.5	15.5
6.7	15.5
6.9	15.5
7.1	15.5
7.3	17.5
7.5	15.5
7.6	15.5
7.8	17.3
7.9	20.0
8.0	17.5
8.1	19.9
8.3	18.9
8.5	15.5
8.6	18.1
8.7	17.3
8.8	17.0
8.9	19.7
9.1	15.5
9.3	15.5
9.5	15.5
9.8	19.1
10.0	15.5
10.2	15.5
10.4	15.5
10.6	15.5
10.8	15.5
11.0	20.0
11.1	15.5

Final apple Size

Diagram IV: Final Size X Size at August 10

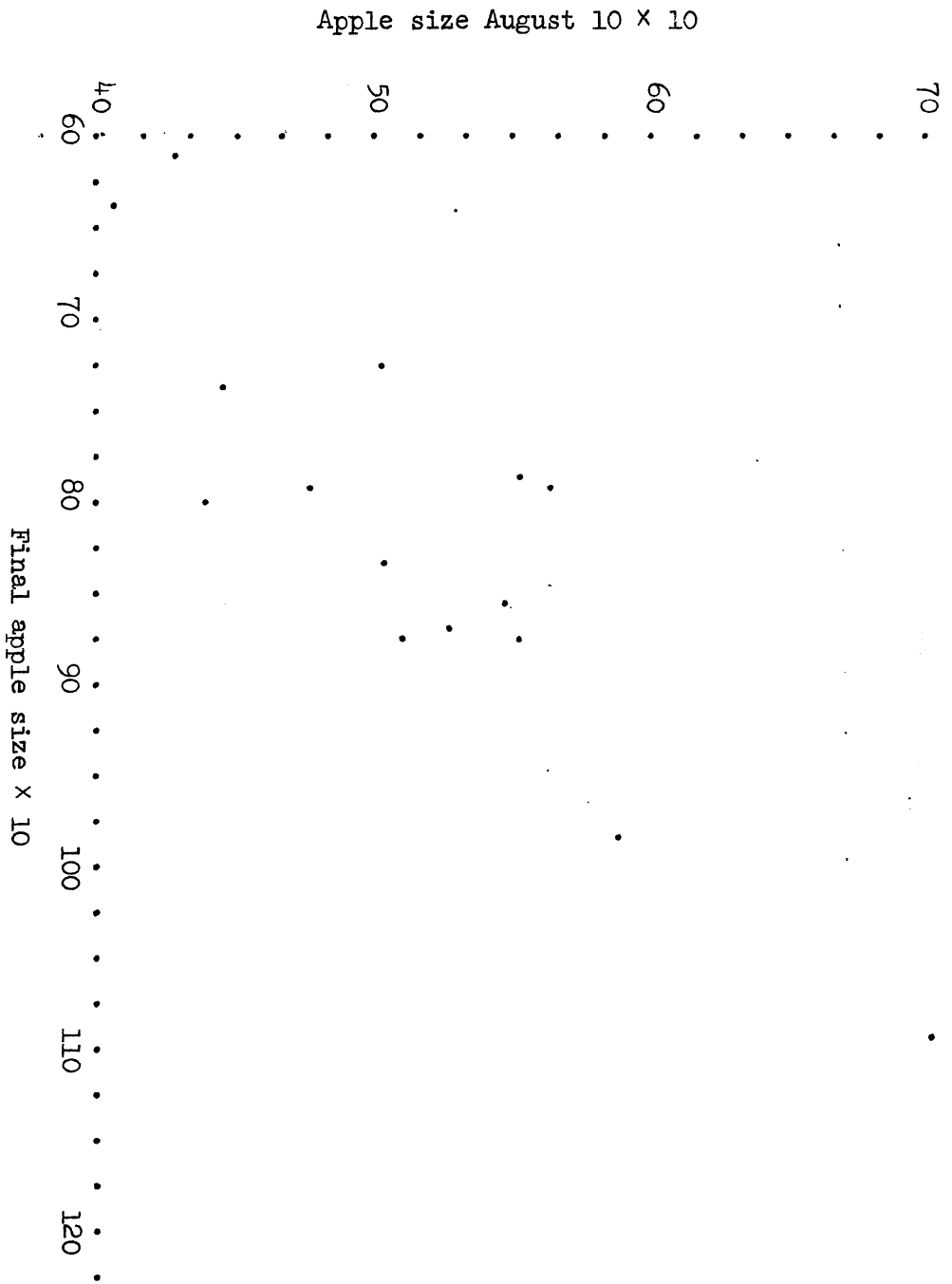


Diagram V: Final Size X Cropload

