الجمهورية العرب المتىة



بمعمدالبخطيط القوى في المرابع

Memo. No. 607

Analysis of Variance

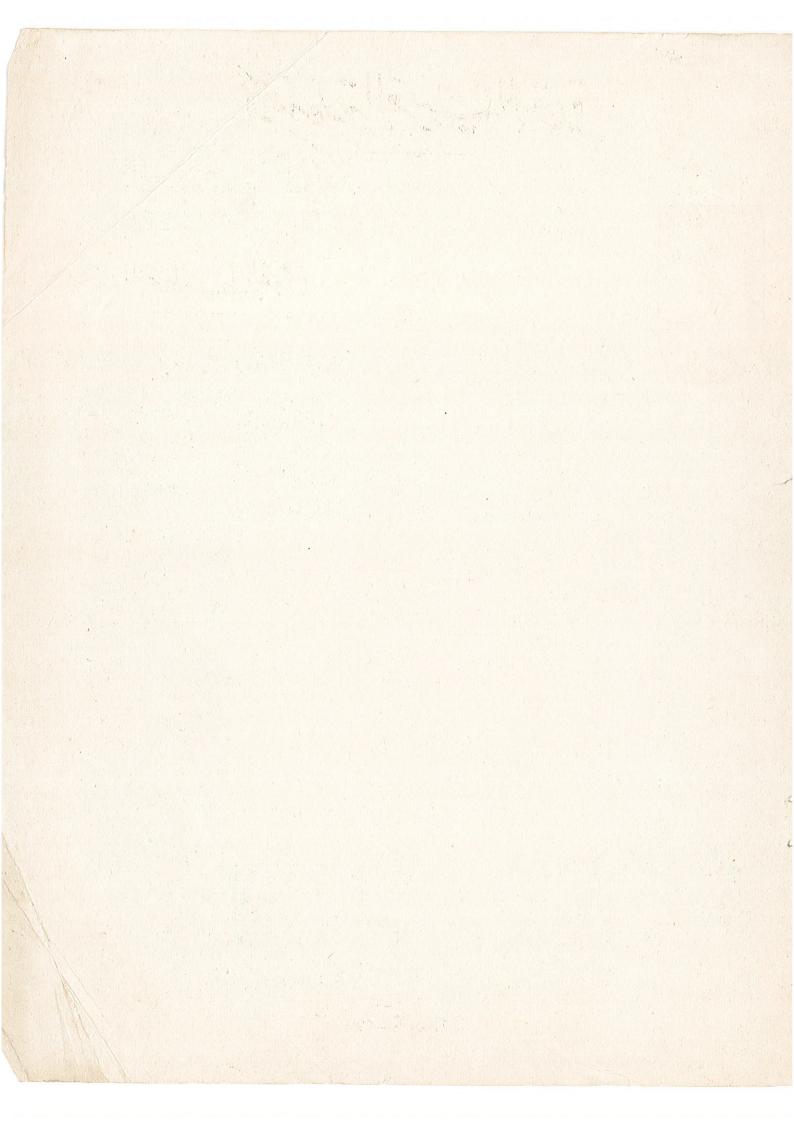
Ву

Dr. Abd El Aziem Anis

December 1965

Operations Research Center.

الق اهدة ٢ شاع محت د مظهر - مالزمالك



Analysis of Variance

I) Def.: The analysis of variance is a statistical technique by which the total variation of the variable being studied can be separated into components that are of experimental interest or importance.

example: as an illustration of the type of problem for which the analysis of variance is useful, consider a gunnary problem experiment in which 4 different brands of shells are to be tested to see whether they are equally satisfactory in quality.

The experiment consists of having 6 different marksmen fire on equal number of rounds with each brand of shell and recording the scores made by each marksman for each brand.

Then scores (the variable) may be arranged in a rectangular array containing 6 rows and 4 columns.

For the purpose of considering other problems, let the scores be displayed in a rectangular array containing a rows and b columns as shown in table.

Table I

x_{21} x_{22} \dots x_{2j} \dots x_{2b} \overline{x}_{2} \dots x_{i1} x_{i2} \dots x_{ij} \dots x_{ib} \overline{x}_{i}			************	THE REAL PROPERTY OF THE PARTY OF THE PARTY.		
x _{il} x _{i2} x _{ij} x _{ib}	\mathbf{x}_{11}	x ₁₂	x _{lj}	0000000	xlp	x ₁ .
x _{al} x _{a2} ····· x _{aj} ······ x _{ab} x̄ _a .	x21	x22 · · · · ·	x _{2j}		x _{2b}	\bar{x}_2 .
x _{al} x _{a2} ····· x _{aj} ······ x _{ab} x̄ _a .						
x _{al} x _{a2} ····· x _{aj} ······ x _{ab} x̄ _a .						
al ac	x _{il}	x _{i2}	x _{ij}		x _{ib}	x _i .
x.1 x.2 x.j x.b x.	^x al	x _{a2}	^x aj		^x ab	x _a .
1 - 2	- -	X	Ī.		x h	ī
	x.1	*.2 -	^.j		^.b	1.

The location of the dot in the index show whether the mean is a row mean or a column mean.

i.e.
$$\bar{x}_i$$
 = $\sum_{r=1}^b x_{ir/b}$

$$\bar{x}_{.j} = \sum_{r=1}^a x_{rj/a}$$

$$\bar{x}_{..} = \frac{1}{ab} \sum_{j=1}^b \sum_{i=1}^a x_{ij}$$

Hence it is obvious that

$$\bar{x}_{\cdot \cdot \cdot} = \frac{1}{a} \sum_{i=1}^{a} \bar{x}_{i \cdot \cdot}$$

$$= \frac{1}{b} \sum_{i=1}^{b} \bar{x}_{\cdot \cdot j}$$

- II) Two well known mathematical models are available for application in experiments of this type:
 - (1) the "Linear hypothosis" model
 - (2) the "components of variance" model.

the essential difference between the two models lies in the assumptions made.

Notice: that x_{ij} is regarded as a set of ab random variables for which the observed values are the values resulting from a single random experiment.

The linear hypothesis model: First assumption: (\propto)

(a) This model assumes that the random variable x_{ij} has a mean \mathcal{M}_{ij} which can be written in the form $\mathcal{M}_{ij} = a_i + b_j + c$

where

$$c = E(\bar{x}),$$

$$a_{i} = E(x_{i} - \bar{x})$$

$$b_{j} = E(x_{i} - \bar{x})$$
(1)

but

$$\sum_{i=1}^{a} (\bar{x}_i - \bar{x}_i) = 0, \quad \sum_{j=1}^{b} (\bar{x}_j - \bar{x}_i) = 0$$

Hence

$$\sum_{i=1}^{a} a_{i} = 0 , \sum_{j=1}^{b} b_{j} = 0$$
 (2)

(1.1) Assumption (1) merely states that the mean of the variable x_{ij} is the sum of a general mean c, a row effect a_i , a column effect b_j .

Application in case of gunnary experiment: This means that if the i-th markman was superior, his mean score would be expected to exceed the mean more for all six marksmen by a + ve quantity a whereas if he were an inferior marksman, a, would be - ve.

Also b_j is a number (+ ve or - ve) which measures the superiority or of the brand j with respect to all brands.

Important criticism: The additive feature of (1) is restrictive. For example, if the rows of table (1) correspond to different amounts of a chemical compound added to the soil, whereas the columns correspond to different quantities of a second chemical compound added, one would not expect the effects of those compounds on crop production to operate independently in this manner.

Second assumption B:

In addition to (1), this model assumes that the variables x_{ij} are independently and normally distributed with the same variances σ^2 . Now if we want to test the hypothesis that the

brands of shells are equally good, this would be expressed in the form that

$$H_0 : b_1 = b_2 = b_3 = \dots b_b$$
In view of (2), this mean
$$H_0 : b_1 = 0 \qquad (j=1, 2, \dots b) \qquad (3)$$

(X) Proceedings of the analysis of variance

Under the foregoing assumptions, we proceed to prove that

$$\sum_{i=1}^{a} \sum_{j=1}^{b} (x_{i,j} - \bar{x})^{2} = \sum_{i} \sum_{j} (\bar{x}_{i,} - \bar{x})^{2} + \sum_{i} \sum_{j} (\bar{x}_{i,j} - \bar{x})^{2} + \sum_{i} \sum_{j} (\bar{x}_{i,j} - \bar{x}_{i,-} \bar{x}_{j,+} \bar{x})^{2} + \sum_{i} \sum_{j} (x_{i,j} - \bar{x}_{i,-} \bar{x}_{j,+} \bar{x})^{2}$$
(4)

The importance of this formulae lies in the fact that its shows that the total variation of the variable x_{ij} could be broken down into three components:

- (1) the first component measuring the variation of row means (ie variation in the marksmen ability).
- (2) the second component measuring the variation of column means (ie the variation in the shell brand's effect)
- (3) the third component measuring the variation in the variables x_{ij} after the row and column effects have be eliminated.

$$(\mathbf{x}_{i,j} - \bar{\mathbf{x}})^2 = \left[(\bar{\mathbf{x}}_{i,} - \bar{\mathbf{x}}) + (\bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}}) + (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_{i,} - \bar{\mathbf{x}}_{i,j} + \bar{\mathbf{x}}) \right]^2$$

$$= (\bar{\mathbf{x}}_{i,} - \bar{\mathbf{x}})^2 + (\bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}}_{i,} - \bar{\mathbf{x}}_{i,j} + \bar{\mathbf{x}}_{i,j})^2$$

$$+ 2(\bar{\mathbf{x}}_{i,} - \bar{\mathbf{x}})(\bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}}_{i,j} + \bar{\mathbf{x}}) + 2() ()$$

$$+ 2(\bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}})^2 = \sum \sum (\bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}})^2 + \sum \sum (\bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}})^2 + \sum \sum (\bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x}}_{i,j} - \bar{\mathbf{x$$

$$\sum_{i=1}^{a} \sum_{j=1}^{b} (\bar{x}_{i} - \bar{x})(\bar{x}_{j} - \bar{x}) = \sum_{i=1}^{a} (\bar{x}_{i} - \bar{x}) \sum_{j=1}^{b} (\bar{x}_{j} - \bar{x}) = 0$$

Since
$$\sum_{j} (\bar{x}_{,j} - \bar{x}) = 0$$
, $\sum_{i} (\bar{x}_{i,} - \bar{x}) = 0$

This applies also to the other two lest sums. Hence (4) is proved.

Application of the F - test:

Before we apply the F-test to these components in (4), we have to convert them into X2 variables.

Let us take variable x ; : It is a linear combination of the basic variables xij which as are assumed to be normal.

Hence x is a normal variable

$$E(\bar{x}_{ij}) = E(\frac{1}{a}\sum_{i=1}^{a}x_{ij})$$

$$= \frac{1}{a} \sum_{i=1}^{a} E(x_{i,j}) = \frac{1}{a} \sum_{i=1}^{a} (a_i + b_j + c)$$

=
$$b_j$$
+ c (since $\sum_{i=1}^a a_i = 0$)

But when H_0 is true $b_j = 0$

Hence

$$E(\bar{x}_i) = c$$

Also since x is a mean of a independent variables

$$V(\bar{x}_{i,j}) = \sigma^2/a$$

Hence all this shows that the variables $\bar{x}_{,j}$ are independently and normally distributed with mean c and variance e^2/a when H_o is true.

It follows then that

$$\sum_{j=1}^{b} (\bar{x}_{.j} - \bar{x})^{2} / \sigma^{2} / a = \sum_{i=1}^{a} \sum_{j=1}^{b} (\bar{x}_{.j} - \bar{x})^{2} / \sigma^{2}$$
 (5)

will possess a χ^2 distribution with (b-1) degrees of Freedom. (on the basis of the assumption H_0).

Now

$$\frac{1}{\sigma^2} \sum_{x_{i,j}} (x_{i,j} - \bar{x})^2 \quad \text{obviously is } \chi^2 \text{ with (ab-1) D.F.}$$

$$\frac{1}{\sigma^2} \sum_{x_{i,j}} (\bar{x}_{i,j} - \bar{x})^2 \quad \text{is } \chi^2 \text{ with (a - 1) D.F.}$$

Hence

$$\sum \sum (x_{i,j} - \bar{x}_{i,0} - \bar{x}_{i,j} + \bar{x})^2$$
 (6)

has a χ^2 distribution with D.F.

$$(ab-1) - [(a-1)+(b-1)] = ab - (a+b) +1 = (a-1)(b-1)$$

Hence if (5) is divided by (b-1), and (6) is divided by (a-1)(b-1), the ratio of the resulting quantities will have an F distribution. It is clear that (5) should be used in testing H_0 because it measures the variation of column mean, and this variation should prove excessively large when H_0 is true as compared to its value when H_0 is true. (6) also should be used because it measures the variation in any other factors and thus should prove useful as a basis for comparison.

Summary of the linear hypothesis F test

If the variables $x_{i,j}$ are independently and normally distributed with means $\mu_{i,j} = a_i + b_j + c$ and variance e^2 , the

hypothesis H_0 : $b_j = 0$ (j=1, 2, ... b) may be tested by using F distribution where

$$F_{i,i,\bar{z}}$$
 (a-1) $\sum (\bar{x}_{i,j} - \bar{x})^2 / \sum (x_{i,j} - \bar{x}_{i,j} - \bar{x}_{i,j} + \bar{x})^2$

and where $V_1 = b-1$, $V_2 = (a-1)(b-1)$.

Notice: The equality of the row means can be tested in the same manner, only F in this case would take a different expression.

ex/ Four plots of lands growing potatoes were divided into 5 subplots each. For each plot 5 treatments were assigned at random to the 5 subplots. The following table was given. Test whether the 5 treayments are equally effective with respect to mean yield.

(Treatment)

		A	В	C	D	E
	1	310	353	366	299	367
	2	284	293	335	264	314
Plots	3	307	306	339	31.1	377
	4	267	308	312	266	342

Test of column mean

$$\sum_{j=1}^{5} (\bar{x}_{,j} - \bar{x})^2 = 3178, \quad \sum_{i=1}^{4} \sum_{j=1}^{5} (\bar{x}_{,j} - \bar{x})^2 = 4x3178$$

$$\sum_{j=1}^{4} (\bar{x}_{i} - \bar{x})^2 = 1286, \quad \sum_{j=1}^{5} \sum_{i=1}^{4} (\bar{x}_{i} - \bar{x})^2 = 5x1286$$

$$\sum_{j=1}^{4} \sum_{j=1}^{5} (x_{i,j} - \bar{x})^2 = 21,530$$

$$\sum_{j=1}^{4} \sum_{j=1}^{5} (x_{i,j} - \bar{x})^2 = 21,530$$

Hence by (4)

$$(x_{ij} - \bar{x}_{i.} - \bar{x}_{.j} + \bar{x})^2 = 21,530 - 4x3178 - 5x1286$$

$$= 2388$$

$$F = \frac{3x4(3178)}{2388} = 16.0 , v_1 = 4 , v_2 = 12$$

The result is significant.

Hence the 5 treatments undoubtedly differ in their effect on yield. row mean:

$$H_0: a_1 = 0$$
 (i = 1, 2,...,a)
 $F = \frac{4 \times 5(1286)}{2388} = 10.8$, $y_1 = 3$, $y_2 = 12$

This result is also significant.

This means that the 4 plots differ in fertility.

The results are usually displayed in a table form

Source of variation	Sum. Sq.	D.F.	M.S.	F
Columns	12,712	4.	3178	16.0
rows	6,430	3	2143	10.8
Remainder	2,388	12	199	
Total	21,530	19		

One way classification:

Now suppose that we are conducting the gunnary experiment to test only whether the different brands of shells are of equal quality.

i.e. Suppose that instead of using six men, we use one man in the test. Clearly here our previous assumptions will reduce to

and we would be testing the hypothesis H_0 : $b_j = o$ (j=1, 2,...b). The analysis of variance formulae would be reduced to

$$\sum_{i} \sum_{j} (x_{ij} - \bar{x})^{2} = \sum_{i} \sum_{j} (\bar{x}_{j} - \bar{x})^{2} + \sum_{i} \sum_{j} (\bar{x}_{ij} - \bar{x}_{j})^{2}$$
 (1)

This could be proved in the following way:

$$(\mathbf{x}_{i,j} - \bar{\mathbf{x}})^2 = (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_{j} + \bar{\mathbf{x}}_{j} - \bar{\mathbf{x}})^2$$

$$= (\bar{\mathbf{x}}_{j} - \bar{\mathbf{x}})^2 + (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_{j})^2 + 2(\bar{\mathbf{x}}_{j} - \bar{\mathbf{x}}) (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_{j})$$

$$\sum_{i} \sum_{j} (\mathbf{x}_{i,j} - \bar{\mathbf{x}})^2 = \sum_{i} \sum_{j} (\bar{\mathbf{x}}_{j} - \bar{\mathbf{x}})^2 + \sum_{i} \sum_{j} (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_{j})^2 + 2\sum_{i} (\mathbf{x}_{j} - \bar{\mathbf{x}}_{j})^2 + 2\sum_{i} (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_{j})^2$$

$$\text{Now } \sum_{i} \sum_{j} (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_{j}) = 0$$

$$\text{Since } \bar{\mathbf{x}}_{j} = \frac{1}{a} \sum_{i=1}^{a} \mathbf{x}_{i,j}$$

Hence formula (1) follows.

It is clear that $\sum (\bar{x}_{ij} - \bar{x})^2/2$ is a χ^2 with (ab-1) D.F. and that $\sum (\bar{x}_{j} - \bar{x})^2/2$ is a χ^2 with (b-1) D.F. Hence $\sum (\bar{x}_{ij} - \bar{x}_{j})^2/2$ is a χ^2 with (ab-b) D.F. and hence $F = \frac{a(b-1)}{b-1} \frac{\sum (\bar{x}_{j} - \bar{x}_{j})^2}{\sum (\bar{x}_{ij} - \bar{x}_{j})^2}$

This could be summarised as follows:

Source of variation	D.F.	Sum of Squares	Mean Square
Between class means	b-1	$\sum \sum (\bar{x}_j - \bar{x})^2$	ΣΣ/b-1
Within classes	b(a-1)	$\sum \sum (x_{ij} - \bar{x}_j)^2$	$\sum \sum b(a-1)$
Total	ab - 1	$\sum \sum (x_{ij} - \bar{x})^2$	

Ex. 35, plots of approximately equal fertility were sown with 7 different variaties of wheat, 5 plots to each variaty, the distribution of variaties among the plots being random. The following table gives the yields of grain in bushels per acre, the 7 columns corresponding to the different variaties. Do the data (fictitious) indicate a significant difference in the yields of the variaties?

13	15	14	14	17	15	16
1.1	11	10	10	1.5	9	12
10	1.3	12	15	14	13	13
16	1.8	13	17	19	1.4	15
12	12	11	10	12	10	11

Answer:

Source of variation	D.F.	S.S.	M.S.	F.
Between variation	6	41.6	6.933	1.1
Within variaties	28	174	6.214	
Total	34	215.6		

$$v_1 = 6, \quad v_2 = 28$$

ex.: The following table given the results obtained from dye trials on each of 5 preparations of Naphthalene Black 12 B made from each of 6 samples of H acid intermadiate.

Yields of Naphthalene Black 12B

Sample of H acid	1	2	3	4	5	6
	1440	1490	1510	1440	1515	1445
Individual yield in	1440	1495	1550	1445	1595	1450
grams of standard	1520	1540	1560	1465	1625	1455
colour	1545	1555	1595	1545	1630	1480
	1580	1560	1605	1595	1635	1520

Does the use of different intermediates gives significantly different yields?

Notice on computation:

Computation could be simplified in many cases:

$$\sum \left[\left(\mathbf{x}_{ij} - \overline{\mathbf{x}} \right)^2 \right] = \sum \left[\mathbf{x}_{ij}^2 - \mathbf{T}^2 \right]$$
 ab (1)

where

$$T = \sum x_{i,j}$$

Also

$$\sum_{j} \sum_{i} (x_{i,j} - \bar{x}_{j})^{2} = \sum_{j} \{x_{i,j}^{2} - T_{j}^{2} / a\}$$

Where T_j is the sum of the values in the j class.

Hence

$$\sum (x_{ij} - \bar{x}_{j})^{2} = \sum_{i} \sum_{j} x_{ij}^{2} - \sum_{j} T_{j}^{2} / a$$
 (2)

Hence

$$\sum (x_{ij} - \bar{x}_{j})^{2} = \sum_{j} T_{j}^{2}/a - T^{2}/ab$$
 (3)

Since the deviations from the means are independent of the choice of origin the results obtained in (1), (2), (3) are unaltered by chance of origin. This would simplify the arithmetic.

ex.: In the first example on fertelises, diminish all the yields by (2)

	1	3	2	2	5	3	4
	-1]	-2	_2	3	3	0
	-2	1	0	3	2	1	1
	4	6	1	5	7	2	3
	0	0	-1	-2	0	-2	-1
Tj	2	9	0	6	17	1	7
×j	0.4	1.8	0	1.2	3.4	0.2	1.4
$\sum_{\hat{1}}$		= \(\sum \) T = 266	•		s 50°4°	1 7 12	= 92.
Hence	by (1)	, (2)	, (3)			
ΣΣ	_(x _{ij}	- x)2	12 2	66 -	50.4 =	215.6	
ΣΣ	_(x _{ij}	- z)2	æ 20	66 =	92 = 1'	74	
		x) ²	= 98	2 - 50.	4 = 41.6		
etc	0						

"Components of Variance" Model

I) Assumptions involved:

This model makes a linearity assumption about the basic variables $\mathbf{x}_{i,j}$ rather than about its mean/ i.j.

In place of (1) (in case of L.H.M.), it is assumed that x_{ij} could be expressed: