Mr. Friedman:

1. This paper used to be CONFIDENTIAL and registered in its previous edition. We sent a letter to the Navy in <u>April</u> requesting this paper, and we still have no answer: this is indicative of the state of chaos existing in their training section.

3. Other than the boor terminology employed, and the plethora of mathematical eyewash that makes a simple subject difficult, this paper is potentially very good, after substantial editing and liberal re-writing.

- Capt 6.

Declassified and approved for release by NSA on 03-10-2014 pursuant to E.O. 13526

REF ID: A66796 From Vision (NCA)

RESTRICTED

# THE INDEX OF COINCIDENCE

#### FOREWORD

The subject of this pamphlet is coincidence.

The student may well ask, "What is coincidence and what applications has it?"

"Coincidence" as the term is used here may be defined as a recurrence of a letter in the same place, or in a corresponding place, as when two texts are lined up one under the other, letter for letter.

This mathematical evaluation assists the cryptanalyst first in preparing his material for attack, and later in the actual attack itself. It assists specifically in answering the following questions.

- 1) How much like random, or how different from random, is this text?
- 2) How similar are these texts?
- 3) How significant is this variation from random?
- 4) How significant is this similarity?

# I SIMPLE MONOGRAPHIC COINCIDENCE

The test of coincidence is the evaluation of the coincidences of letters, or of digraphs, etc., between two or more messages, or within the same message.

The coincidence or "pairing" test may be consolidated into one final number or "statistic". That statistic is called the "index of coincidence" and is defined as the ratio of the actual coincidences to the coincidences to be expected from chance (coincidences in random text). For English text the expected I.C. is 1.75. For most European languages the expected I.C. is about 2.00. For random text the expected I.C. is 1.00.

Assume two pages of cipher text based on a complex cipher which will give a "flat" frequency table for the entire message. Select a letter at random (say the 3rd) from one page and another from the other page (say the 3rd also).

There is 1 chance in 26 of the first letter's being an "A". There is 1 chance in 26 of the second letter's being an "A" There is 1 chance in 676 of both letter's being "A" There is also 1 chance in 676 of both letter's being "B"

Therefore, the chances of both letters being the same letter (in a chance selection of cipher text) are:

26 chances in 676, or 1 chance in 26, or 3.8462%.

If we select many pairs of cipher letters, the average number of identical letters to be expected "in the long run" will be 3.846% (or 1/26) of the total number of possible coincidences. We call this number the "Expected Coincidence due to Chance" (random text).

With English text it is different. Take two pages of English text. Make a chance selection from each page.

There are about <u>130</u> chances in <u>1,000</u> of the first letter's being an "E" There are about <u>130</u> chances in <u>1,000</u> of the second letter's being an "E" There are about <u>16,900</u> chances in <u>1,000,000</u> of both letters' being an "E" Likewise, there are <u>8,464</u> chances in <u>1,000,000</u> of both being "T" <u>6,400</u> chances in <u>1,000,000</u> of both being "N", etc.

(See table following).

Table

Any letter	1,000	1,000	66,930
Z	1	1	· <u>1</u>
Q X	5	5	25
	2	2	. 4
J	2 、	.2	4
ĸ	4	4	16
Ŵ	16	16	256
B` V	15	15	225
B`	īø	ĩø	100
Ŷ	19	19	361
Ĝ	16	16	256
P F	. 28	28	784
p	27	27	729
ົບ	26	26	676
С М	25	25	625
Č	31	- 31	1,296 961
L	36	36	
H	34	34	1,156
ת	42	42	1,764
T	92	92	· 8,464
S	61	61	5,776 3,729
R	76 -	76	6,100
Ň ·	80	8Ø	6 100
I	73	. 73	5,329
Ă	73	73	5,625 5,476
· E O	13Ø 75	130 75	16,900
	-ing this itr.	ing this ltr.	this letter
iext/	ing this ltr.		letters' being
(Telegraphic Text)	letter's be-	letter's be-	
fext Letter	Chances in 1,000 of 1st	Chances in 1,000 of 2nd	Chances in 1,000,000 of both

Finally there are 66,930 chances in 1,000,000 (the sum of the chances for the individual letters) of both letters' being the same plain text letter in a chance selection. Therefore, if we select many pairs of plain text letters, the average number of identical letters to be expected "in the long run" will be 6.693% (about 1/15) of the total number of Possible Coincidences.

We may call this number the Expected Coincidences in English Text

In actual practice we are concerned with the coincidences between our two texts, or within our alphabet, etc. The tally or count of these coincidences we call the Actual Coincidences.

To permit comparisons between results obtained from texts of varying amounts, it is most convenient to convert to an index number. We call this the Index of Coincidence and use the abbreviation  $I.\overline{C. \text{ or } \iota}$ .

By definition / = Actual Coincidences Expected Coincidences due to Chance.

The expected I.C. for English (or mono-alphabetical cipher text) is: <u>.06693</u> = 1.75, approximately.

The actual I.C. of unknown cipher text may take almost any value but in practice the range will generally extend from about .80 to about 2.00 (simple monographic index of coincidence).

The value of the index of coincidence for a given English text will depend on the actual distribution of letters in that text. Repetitions in short texts will increase the index of coincidence. Unrelated text (that is, text with few repetitions) will give an I.C. approaching the theoretical 1.75. As the expected number of chance coincidences is based on a flat frequency (where each cipher letter is ultimately used the same number of times) any cipher text that differs radically from such frequency distribution will have a correspondingly higher I.C. This is especially noticeable in short cipher texts where the frequency table has not had an opportunity to "flatten out".

The mono-graphic I.C. of English naval text will increase with small amounts of text to 1.80 - 2.00 (as compared with the theoretical 1.75) and small amounts of random text will give I.C.'s of 1.10 - 1.20 (as compared with the theoretical 1.00). The amount of excess attributable to the sample size will be discussed later, under "standard deviation".

For most European languages the expected I.C. is higher than in English, due to the more irregular letter distribution of their normal alphabets, namely:

Language	<u>1.C.</u>
Random text	1.Ø
English	1.7
Italian	1.9
Spanish	2.Ø
French	2.0
German	2.0

#### II POLYGRAPHIC COINCIDENCE

In addition to the simple monographic index of coincidence (i), there are occasions when the digraphic index of coincidence  $(i_2)$ , trigraphic I.C.  $(i_3)$ , tetragraphic I.C.  $(i_4)$ , pentagraphic I.C.  $(i_5)$ , etc., can be used to advantage. They are derived from the normal digraphic (trigraphic, etc) frequency tables in the manner indicated in paragraphs 3 to 7.

Expected values for these simple polygraphic indices of coincidence are as follows:

Language	L	42	63	14	15
Random text	1.00	1.00	1.00	1.00	1.00
Englisa	1.75	4.75	27,89*	х	х
Italian	1,92	5.68	X	Х	X
Spanish	2.02	6.29	Х	X	х
French	2.02	6.29	X	X	Х
German	1,98	6.57	X	Х	Х

Notes: X = Not computed.

2

=(Computed from the only known trigraphic table. (The correct index might vary widely from this estimate)

In practice the actual polygraphic I.C.'s will usually run higher than their theoretical values, and a repeated word or two in short texts will made them sky rocket. As typical examples, we have taken the plain text of four problems in the elementary and secondary courses and computed the various I.C.'s (from  $\iota$ , to  $\iota_s$ , that is the monographic, digraphic, trigraphic, tetragraphic and pentagraphic indices of coincidence).

Text	L	L2	63	14	٤5.
Expected random	1.00	1.ØØ	1.ØØ	1.ØØ	1.ØØ
Expected plain	1.75	4.75	27.89	?	?
Problem No. 1	1.8Ø	5.23	29.11	427.	724Ø.
Problem No. 2	2.ØØ		66.Ø4	1062.	149ØØ.
Problem No. 3	1.91	5.6Ø	42.Ø4	666.	12Ø7Ø.
Problem No. 4		4.9Ø	31.7Ø	456.	919Ø.

#### III THEORETICAL RECAPITULATION

Suppose we have a language for which we know the overall proportions of the letters are  $P_1$ ,  $P_2$ , - - -,  $P_c$ .

(1)

$$\sum_{i} p_{i} = I.$$

Suppose further that we have two pieces of text from this language and line them up one above the other, and then count coincident letters. What is the expected number?

At a particular place the probability of a coincidence involving the i th letter is  $Pi^2$ . There fore, the C cases being mutually exclusive, the probability of an incidence is

$$\sum_{i} \rho_{i}^{2} \qquad (2)$$

If the length of overlap is N, then the expected number of incidences is

N 5 2

If the text is such that 
$$p_i' = p_j = \frac{1}{2}$$
, we will refer to it as "flat", "random". The probability of an incidence is  $\sum_{i=1}^{2} \frac{1}{2} e_i = \frac{1}{2}$ , and the

expected number is N  $\frac{1}{C}$ . The ratio of the number found in a comparison to that expected is called the

The ratio of the number found in a comparison to that expected is called the "index of coincidence", .

$$\iota = \frac{g}{N_c} = \frac{cg}{N} . \tag{3}$$

The expected value  $\gamma$  of  $\iota$  for our language is given by taking the expected value of  $g = N \sum_{i}^{c} \rho_{i}^{2}$  over the expected value for flat text or  $\gamma = \frac{N \sum_{i}^{c} \rho_{i}^{2}}{N_{\Delta}} = C \sum_{i}^{c} \rho_{i}^{2}$  (4)

Notice that the expected value of the I.C. for flat text is 1.

#### IV PRACTICAL APPLICATIONS

#### (A) TO DETERMINE WHETHER TWO MESSAGES ARE IN THE SAME KEY

During U.S. Fleet Problem V (1925) the Battle Fleet used a cipher of their own design. A total of 13 messages in this cipher were submitted to the Code and Signal Section for attack. Although a different indicator was used in each case, it was suspected that some of the messages might be in the same key. Two messages in one key (example No. 1) and two more in another key (example No. 2) were discovered. (The messages were eventually solved).

Each message was "lined up" with each other message and the coincidences were noted. (See examples No. 1 and No. 2).

Example No. 1

or

						G G													
						M O													
						P D													
						J <u>J</u>													
						D Q													
~	-	-	-	_	·	-	-	-	-										

#### Q Q U I I V H B K Q

Coincident letters are underscored. 12 coincidences in 140 pairs of letters. -

Simple monographic coincidence.

Expected	N	3	<u>14ø</u>	2	5.4	Where N = number of units examined,
Coincidences	C		26		•••	C = number of cells for single letter
						examination = 26.

IC =  $\frac{12.0}{5.4}$  = 2.2 (Messages are in same key).

There is one repeated trigraph, GPZ, in the messages under examination. This coincidence indicates that the keys correspond at that point, but does not necessarily indicate that the keys correspond throughout the message. To prove coincidence of the keys throughout the 'two messages, we must have our coincidences spread through the messages in question. (As they were in the above example). Likewise, digraphic and trigraphic coincidences may be compared and evaluated to an index of coincidence.

For example, in the above messages, 2 coincident digraphs were found (GP and PZ) (also one coincident trigraph). In this message there were 139 digraphs and 138 trigraphs in alignment with possibilities of coincidence.

$$\frac{N}{676} = \frac{139}{676} = .206$$

digraphs were to be expected from chance. Two were found, giving an IC =  $\frac{2}{.206}$  = 9.70.

This value, far above the normal 4.75 index of coincidence, does not necessarily indicate the messages are in the same key. All we really know is that the two keys are identical in the second group. The extremely high I.C., 9.70, is due entirely to the small amount of text involved in this example. As the amount of text decreases, the variation of the I.C. from the expected will become more pronounced, until at times it is possible <u>that small</u> amounts of text may give entirely false indications. This effect will be discussed more fully under "standard deviation".

Example No. 2

														J W																
P	E	K	D	A	P	N	U	T	រ	D	D	U	W	<b>Q</b>	Q	I	т	N	P	e	X	t	G	H	T	K	D	C	L	
J	S	D	N	D	U	J	K	L	ប	Z	J	C	S	<b>Q</b>	H	I	0	Z	H	U	K	E	G	X	D	E	P	W	T	
Z	R	P	W	A	M	M	T	I	Q	J	F	P	K	F	D	C	0	V	P	D	Ū	С	H	Z	W	X	M	G	E	
P	R	F	R	X	P	I	Q	V	A	F	Q	R	P	F	E	A	Q	P	Q	E	W	0	Y	E	G	X	R	O	Y	
														A 0																
N	M	K	0	C	Z	D	I	X	B	Y	A	Ľ	V	S	P	G	P	X	F	N	U	E	F	X	N	W	D	F	L	
D	T	E	P	P	D	Z	G	T	J	O	C	V	H	T	P	M	R	P	H	T	Y	C	Q	X	L	L	V	R	N	
														R. W																•
Q	A	W	S	Q	j	N	G	Р	E	O	A	M	D	N	Z	A	A	P	F	Á	Z	Y	0	T	N	Q	W	V	X	
U	M	N	G	C	F	J	F	L	Q	R	M	E	D	F	O	Z	N	S	P	E	C	D	S	L	I	Z	Q	Y	A	
								Y B					U -	0 -	s -	Q	D -	N -	0 -	B -	A -	s -	H -	H -	c ,	т -	D -	D ~	U -	

14 coincidences in 220 pairs of letters.

 $\underline{N} = \underline{220} = 8.46$  coincidences expected. C 26

IC =  $\frac{14}{8.46}$  = 1.66 (almost normal for English).

These two messages <u>probably</u> are in the same key (and actually proved to be). Note that there are no repeated digraphs or trigraphs. Note also that coincidences are well spread out.

- 6 -

# (B) TO DETERMINE WHETHER TWO MESSAGES "OVERLAP" IN THE SAME RUNNING KEY

Copy each message on a single line, omitting all spacing. "Line up" the messages and note coincidences. Then shift one message one place to the right and note the coincidences. Repeat this process to the end.

If the index at any point is 1.75 or higher, for mono-graphic examination, the position and the fact of the "overlap" is probable. In this application the digraphic and trigraphic indices are useful adjuncts to the monographic index.

For some purposes the fundamental unit may be taken to be a set of letters, as digraphs, trigraphs, etc. Suppose we are interested in digraphic coincidences. Then the digraphic I.C.,  $\iota_{\ell}$ , will be calculated as above, noting that the size C of the alphabet is larger this time.

### V RELATIONS AMONG THESE STATISTICS

The monographic and digraphic I.C.'s are not independent. For if the probabilities of the various letters are  $P_1$ ,  $P_2$ , - - -,  $P_c$ , then the probability  $\gamma_{ij}$  of the  $ij^{th}$  digraph is  $\rho_i \rho_j$ , ignoring the cohesion of the language, and for the moment treating it like newspaper which had been cut into little pieces, one letter to a piece, and then shuffled and arranged in a line. Using this estimage of  $\gamma_{ij}$  we get the I.C.

$$\begin{split} \gamma_{2} &= C^{2} \sum_{i}^{C} \sum_{j=i}^{C} \gamma_{ij}^{2} = C^{2} \sum_{i}^{C} \sum_{j}^{C} \rho_{i}^{2} \rho_{j}^{2} \\ &= C^{2} \sum_{i}^{C} \rho_{i}^{2} \sum_{j}^{C} \rho_{j}^{2} = \left( C \sum_{i}^{C} \rho_{i}^{2} = \gamma^{2} \right) \end{split}$$
(5)

It is true however that language has cohesion, and that each letter affects the probability of occurrence of others in its vicinity. Usually then  $\ell^2$  is in excess of the estimate  $\ell^2$  above. We will sometimes calculate the ratio

 $\mathcal{L}_{l_0}^{\ell_0} = \chi_{p}$  and call it the digraphic "index of cohesion".

Estimates can be made in the same way for higher I.C.s. One can show that  $\gamma_{j} = \iota_{j-1} \iota$  or

 $\gamma_{i}^{\prime} = \iota_{i-j} \ \iota_{j} \tag{7}$ 

In these equations the right members are thought of as quantities already computed; while the  $\mathcal{N}'_{\mathcal{S}}$  on the left are estimates or predictions of quantities which can be computed from the definition (4).

An application of these relations occurs in the study of fractionating\* systems, where as a preliminary to enciphering the text is expressed as a product of two components, and each component is enciphered separately, and then the cipher text is recombined to ordinary letters. For instance, each of 25 letters may be represented by a two digit number, where the first digit comes from 0, 1, 2, 3, 4, and the second from 5, 6, 7, 8, 9. The argument for these proceeds as for digraphs, and the I.C. of the combined text is the product of those of the fractions. However, here any significant deviation from this estimate is not adequately described as cohesion, but must be due to the dependence of the, two fraction streams.

### VI THE ROUGHNESS OF A SINGLE SAMPLE

We have introduced the I.C. as a measure of the match between two pieces of text. We can extend this idea now to a measure of the roughness of a single sample. Suppose we have a piece of text which we duplicate on two slips of paper and then place them one under the other for the purpose of counting coincidences. There will be one position of total coincidence, which we will rule out. If we compute the I.C. for all other positions, we will have what we call the "index of coincidence of a single sample.

\*Consult Gaines "Elementary Cryptanalysis," Chapter XXII

L

(6)

If there are M letters in our sample, we will have looked in N = 1/2M(M-1) distinct places.

If the text is flat we would expect  $\frac{N}{C} = \frac{M(M-1)}{2C}$  coincidences.

If the text is not flat but have proportions  $P_1$ ,  $P_2$ , ---,  $P_c$  of letters, then there are  $f_i = P_t M$  occurences of the i th letter. In the course of our counting we will (10) compare every letter with every other, so that the ith letter will give rise to 1/2 (10)  $f_t$  ( $f_t$ -1) coincidences, or

$$\sum_{l}^{c} \frac{1}{2} f_{l} (f_{l} - l)$$
 (11) in all.

(8)

(9)

Comparing this with the expected in flat text we get the I.C.

$$\vartheta = \underbrace{\sum_{i} f_{2}}_{i} f_{i} (f_{i} - 1)$$

$$\frac{\vartheta}{2c} M (M - 1)$$

$$\delta = \underbrace{c \sum_{i} f_{i} (f_{i} - 1)}_{M (M - 1)}$$
(12)

In theoretical context (12) is more useful than (13), but (13) is a little simpler for computational purposes.

Notice that this formula is different from (4). This is because of the omission of the perfect hit. If M is large enough then  $f_l -1$  can be replaced by  $f_l = P_1 M$  and M-1 by M, so that

$$\delta = \frac{c}{\frac{c}{M^2}} \sum_{i=1}^{c} \frac{f_i^2}{i} = c \sum_{i=1}^{c} p_i^2 = \gamma$$
(14)

15 an asymptotic expression for  $\delta$ Notice that  $C \ge \gamma \ge l$ 

We see that

١.

$$\delta \frac{M-I}{M} = \gamma - \frac{C}{M}$$
or  $\delta (M-I) = \gamma M - C$ 

$$\delta = \frac{\gamma M - C}{M} \quad \text{(15)}$$
(16)

 $\delta = \frac{M}{M-1}$  or  $\gamma = \frac{M}{M}$ 

The error in using  $\gamma$  (usually more convenient) in place of

$$\delta \quad is \qquad \gamma - \delta = \frac{c - \gamma}{M - i} = \frac{c - \delta}{M}$$
 (17)

This error is always positive, that is,  $\gamma$  is an over-estimate of  $\delta$ . The error is smaller for larger values of  $\gamma$ , or larger values of M.

Notice that  $\gamma$  is a measure of the shape of the distribution only, and is independent of the sample size, as is

$$\gamma = \frac{c \ge f_{1}^{2}}{M^{2}} = \frac{c \ge p_{1}^{2} M_{1}^{2}}{M^{2}} = c \ge p_{1}^{2}$$
(18)

cince  $f_{l} = P_{j}M$ .

But  $\hat{O}$  does depend on the sample size M, which is a desirable characteristic, since random roughness is usually present in small samples. For smaller samples

$$\delta' = \gamma - \frac{c - \gamma}{M - 1} \tag{19}$$

is seen to be smaller, thus automatically compensating to some degree for small sample errors. We will usually measure the roughness of single samples by  $\sigma$ , using  $\gamma$  as an asymptotic approximation.

#### VII EXAMINATION OF CIPHER ALPHABETS AND CIPHER TEXTS

The indices of coincidences discussed in the previous paragraphs may be used in analyzing the internal structure of a cipher alphabet. A message of 173 letters has a frequency table as given below:

A	В	С	Ð	Ł	F	G	Н	I	J	К	L	М	N	0	р	Q	R	S	Т	U	v	W	х	- Y	Z
5	3	ิด	14	2	2	10	22	6	1	8	13	1	ø	14	16	ø	13	2	13	Ø	7	19	1	2	1

We make a tally count. i.e., the number of times a letter occurs with the same frequency, thus:

<u>Tally</u> f	Number of <u>Tallies</u> n	$\frac{f(f-1)}{2}$	$\frac{n(x) f(f-1)}{2}$
Ø	4	A	0
<u>}</u>	3	ต	, <u>(</u>
2	- 4	1	4
3	1	3	<u>?</u>
4.	1	6	6
6	2	15 .	· 30
7	. 1	· 21	. 21
8	1	·28	28 •
10	2	, 45	90
13	3	78	234
14	2	91	182
19	1	171	171
22	1	231	_231
			Crincidences 1000

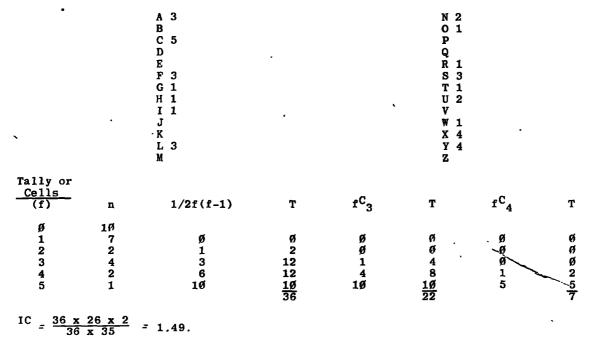
We find that there are 1000 coincidences. Although we can not count the coincidences in the examination of a single cipher text, we can evaluate the various frequency counts into actual coincidences. Having the actual coincidences from the table (1000), we obtain the 1.C. by formula (13)

$$\delta = \frac{c \sum f_i(f_i - 1)}{M(M - 1)}$$

Substituting.  $\delta = \frac{1000 \text{ x} 52}{173 \text{ x} 172} = 1.75$ 

The index of coincidence indicates that a monoalphabetic substitution was employed.

As a second example to show the results obtained from small texts, we calculate as follows from a frequency count of 36 letters assumed to be monoalphabetic.



The alphabet in question was actually a monoalphabetic substitution. With a small amount of text, the simple index is somewhat indeterminate. Using the triple and quadruple indices, the results are even more so, and at times may give even false indications. It is again emphasized that sufficient text must be used to give positive indications.

• •

As another elementary example of the application of the index of coincidence to the internal examination of a cipher text, we have, for example, a 5-letter repetition at an interval of 85. Is the cipher a polyalphabet cipher of 5 or 17 alphabets? By means of internal examination with the index of coincidence we can determine what type of cipher we have. Make a frequency count of the cipher alphabets assuming 5 and then 17 alphabets.

Calculate the index of coincidence in each case for one or more alphabets. The indices of higher value will indicate which assumption is correct. If neither assumption shows positive results (an index around 1.7) we may have a progressive cipher, running key cipher, auto key cipher, or cipher of even more complex nature.

Another elementary application is as follows. We have a cipher message which has been intercepted. The I.C. is computed and found to be  $O^{-1.79}$ .

This is so rough as to resemble plain text. A simple substitution has the property of leaving O' unchanged, and so has a transposition. Multi-alphabet substitutions lower the I.C.. So we are reduced to three hypotheses, that our sample is either transposed plain text, a simple substitution, or both substitution and transposition.

The digraphic I.C. is computed,  $\sigma_2 = 4.85$ . Remember that its expected value is  $3.05 = \sigma^2$  in view of the known roughness. Therefore the index of coherence is  $K_2 = \frac{4.85}{3.05} = 1.56$ 

Since a transposition destroys coherence we can assert that no transposition is involved. Multigraphic Indices of coincidence are preserved by a simple substitution.

### VIII THE STANDARD DEVIATION

We have already several times referred to the fact that these statistics are useful only if the sample is large enough. To get an idea as to whether this is the case or not we measure our results in terms of a standard deviation, "sigma". One standard deviation is roughly one half the width of a band which when placed about the average will include two-thirds of the data. It is a measure of the dispersion. If sigma is large the data is spread out wide, and if it is small the numbers are close together. In a binomial distribution the standard deviation is  $O = \sqrt{N} \times \frac{\rho q}{\rho}$  where N is the number of observations and  $\rho$  and q are the probabilities of success and failure. -

To estimate the significance of  $\mathcal{O}$ , we refer to (12), where the denominator is the expected number of incidences in flat text, and the numerator is the number found. Assuming a binomial distribution of the incidences we find the variance  $\mathcal{O}^{-2}=N_{pq}=\frac{M(M-1)(C-1)}{2C^2}$  (20)

If s is the "sigmage" or deviation of the number found divided by  $\sigma$  we have

$$s = \frac{\frac{1}{2} \sum_{i} f_{i} (f_{i} - 1) - \underline{M(M - 1)}_{2C}}{\sqrt{\frac{M(M - 1)(C - 1)}{2C^{2}}}}$$
(21)

$$\frac{c \sum f_{i}(f_{i}-l) - M(M-l)}{\sqrt{2(C-l)(M)(M-l)}} = \frac{c \sum f_{i}(f_{i}-l)}{\frac{M(M-l)}{\sqrt{2(C-l)}}} = \frac{\partial -l}{\sqrt{2(C-l)(M-l)}}$$
(22)

$$=\frac{\delta - 1}{\sqrt{2(C-1)}} \sqrt{M(M-1)} = \frac{\delta - 1}{\sqrt{2(C-1)}} M.$$

For M>51 error is less than 1%.

Notice that the sigmage is a linear function of the sample size M, and also linear with the "bulge"  $\mathcal{O}_{-}/$ . The denominator is relatively unimportant to the estimation of S except in shifting from code to cipher, when C can change from 500,000 to as small as 10.

The bulge  $\mathcal{O}_{-}/$  is a quantity which will recur frequently.

Formula (23) does not apply to the iota I.C. For that we have the expected number g and  $\frac{N}{c}$  found. Then  $\sigma^{e} = \frac{N}{c} \left( 1 - \frac{1}{c} \right)$  and the sigmage is

$$S = \frac{g - \frac{N_{c}}{N}}{\sqrt{\frac{N_{c}}{1 - \frac{1}{c}}}} = \frac{g_{c}}{\sqrt{\frac{G - 1}{N}}} = \frac{1 - 1}{\sqrt{\frac{G - 1}{N}}}$$
(24)

In this case the signage is linear with the bulge l - / , but varies only as the square root of the sample size.

The significance of s is given in the following table, which lists the probability of getting s or a larger result from chance.

(23)

s	prob.	S	prob.
	-		
.1	.4602	2.9	.0019
.2	.4207	З.Ø	ØØ135 = 1/8ØØ
.3	.3821	3.2	, ØØØ69
.3 .4 .5 .6	.3446	3.4	.00034 - 1/3000
.5	.3085	3.6	.00016
.6	.2743	3.8	. ØØØØ7
.7	.2420	4.0	.00003 = 1/33.000
.8	.2119	· 4.1Ø1	.0000206
.9	.1841	4.200	.0000134
1.0	.1587	4.299	.0000086
1.1	.1357	4.398	.0000055
1.2	.1151	4.497	.0000035 = 1/300,000
1.3	.0968	4.596	.0000022
1.4	.0808	4.695	.0000014 = 1/711,000
î.5	.0668	4.794	.0000008 - 1/1.250,000
1.6	.0548	4.907	.00000046164
1.7	.0446	5.006	.00000027741
1.8	.0359	5.105	.00000016513
1.9	.0287	5,209	.00000009736
2.Ø	.0228	5.303	.00000005686 = 1/18 million
2.1	.0179	5.402	.00000003290 =1/30 million
2.2	.0139	5.501	.666666961856
2.3	.0107	5.600	.0000001070
2.4	.0082	5.798	.00000000335
2.5	.øø62	6.Ø8Ø	.00000000060
2.6	.0047	6.5Ø3	.00000000004
2.7	.0035	6.785	.00000000001 = 1/100,000 million
2.8	.0026		·····

#### EXAMPLE

Suppose an unknown cipher in four digit groups is being investigated and the question is whether it is reenciphered\* or not. If it is the text can be expected to resemble random more than if it is not. We make a frequency count on the 10,000 groups and then determine how likely such a distribution is by chance. If it is not likely we must seek an explanation.

If 560 groups are counted and  $\sum f_i(f_i-I) = 76$  then  $\mathcal{O} = 2.4I$  This gives a sigmage s = 56. The table shows that we would have to repeat this procedure about 120 million times on random material to get a like result. We can say that this result is unlikely by chance and that an explanation is called for. The obvious one is that there is no reencipherment, or that it is very feeble. Tests to check this further hypothesis can be quickly devised.

\*See Gaines "Elementary Cryptanalysis" page 2. 👘

1

# IX TO DETERMINE WHETHER TWO ALPHABETS ARE IDENTICAL ALPHABETS

Assume that a complex cipher using secondary alphabets has been analyzed and reduced to  $5\theta$  alphabets. There are only 26 possible secondary alphabets, so some of these  $5\theta$  alphabets can be combined. Visual inspection is too inaccurate to be trusted, except within an abnormally large amount of text.

Four sample alphabets are given in the table following, A and B.

### Table "A" -- Frequency Tables

No.	1		No	<u>). 2</u> .	<u>No.</u>	<u>3</u>	<u>No.</u>	4
Δ			А		А	2	A	2
B	2		B		B	2	в	
č	2 1		· c		Ē	2 2 1	A B C D	1
ň	-		n	1	- D	-	D	1 1
Ē	3		Ē	1 3	Ē	2	Е	
F	3 1		A B C D E F	•	A B C D E F G H	-	E F G H	1
Ĝ	-		Ğ		Ē	2 .	G	
ਸ	2		H		, Н		н	
A B C D E F G H I J K	2 2 1		Ï		ī	3	I	2
Ĵ	1		I J	1	J	-	J	
ĸ	•		ĸ	2	ĸ	•	к	
L			L	1 2 1	L		L	
M	2		M	-	. M	1	М	
M N	2 3		M N	2	N N	1 3		1
	-			_	0		N P Q R S	1 2 1
P			P	1	P		Р	1
ā			ā	ī			Q	
Ř			R	1 1	R		R	
S			S		Q R S		S	3
O P Q R S T U V			O P Q R S T	1	Т	•	Т	
`Ū			U		U	· .	U	
v		1	U V	2	. <b>v</b>		v	
W			W		h.	1	<b>W</b> .	3
х			Х	2 2	Х		W X Y	3 2 1
X Y	2		W X Y	2	X Y	1	Y	1
Z	1		Z		Z	2	Z	
	0.0		_	20		2Ø		2Ø
	2Ø			20		2 P		~~

### Table "B" -- Repeated Letters

<u>No.1-&amp;-</u>	<u>2 No</u>	<u>.1-&amp;-3</u>	<u>No.1-&amp;-4</u>	<u>No.2-&amp;-3</u>	<u>No.2-&amp;-4</u>	<u>No.3-&amp;-4</u>
A	A		A B C 1	A	A	A 4
B	A B C	4	B	B	B	B
A B C D E 9		1	C 1	С	С	C 1
D	D		D	D E 6	D 1	D
E 9	Ē	6	E F 1	E 6	E	E
F	F		F 1	F	F	F
G	G		G	G	G	G
Н	н		н	н	н	Н
I	I	6	I 4 •	I	I	16
J	J		J	J	J	J
K	K		K	К	K	K
L	L		L	L	L,	L
М	М	2	М	М	M	M
N 6	N	9	N 3	N 6	N 2	N 3
0	0		0	0	0	0 '
Р	Р		Р.	Р	P 1	Р
Q R	Q		Q	Q	Q.	Q
R	R		R	R	R	R
S	· S		S	S	S	S
Т	Т		Т	T	Т	Т
U	ប		U	U	ប	υ.
v ·	v		V	"V	V	V
W	W		W .	w	W.	W 3 '
х	х		X	X	X 4	Х
<u>Y 4</u>	Y		<u>Y 2</u>	<u>Y .2</u>	<u>Y 2</u>	Y 1
		32	11	14	1Ø	18 .
	oin-					
c	idences			•		
	aa i aa i da					

Chance coincidences =  $\frac{400}{26}$  = 15.4

- 13 -

Line up two alphabets at a time and cross multiply the repetitions for each letter (see table "B"). "E" occurs 3 times in No. 1 alphabet and 3 times in No. 2. There are 9 pairs of "E's" in No. 1 and No. 2. Add the coincidences noted for a pair of alphabets.

There are 20 letters in each alphabet which gives 400 pairs of letters to deal with. Chance would give 400/26, or 15.4 repeated letters in two alphabets. The index of coincidence is the sum of the actual coincidences divided by 15.4.

The induces are as follows:

 No. 1-&-2
 1.23

 No. 1-&-3
 2.08
 (above the normal 1.75)

 No. 1-&-4
 .71
 (alphabets No.1-&-3 must be identical)

 No. 2-&-3
 .91

 No. 2-&-4
 .65

### X THE CROSS I. C.

Suppose we have two stretches of text of which the distributions are given by  $P_1$ ,  $P_2$ , ---,  $P_c$  and  $q_1$ ,  $q_2$ , ---,  $q_c$ 

$$\sum_{i}^{C} p_{i} = 1, \quad \sum_{i}^{C} q_{i} = 1.$$

If these are placed one above the other the probability of an incidence at any one position is  $\mathcal{C}$ 

$$\sum_{l=1}^{p_i} p_i q$$

and the I.C. is

$$\xi = c \sum \rho_i q_i$$

We can show that the expected value of  $\zeta$  is 1. Suppose that we apply a permutation to the q's and recompute  $\zeta$ . If we do this for all C: permutations and average them we get

$$E(\xi) = \frac{1}{C}\sum_{\substack{A|l\\permutations}} \xi = \frac{1}{C}\sum_{\substack{A|l\\permutations}} c \sum_{i}^{C} p_{i} q_{i}'$$

where  $q_{l}'$  is one of the q's, depending on the permutation. Each q comes into a given position (C-1)! times. Thus

$$\sum_{\substack{A_{1} \\ p \in mutations}} q_{L} = (C - I)! \text{ and}$$

$$E(\xi) = \xi_{L} \sum_{L=I}^{C} p_{L} \sum_{\substack{A_{1} \\ P \in mutations}} q_{L}' = \frac{C(C - I)!}{C''} \sum_{I} p_{I} = I$$

Notice that this result is independent of the roughness of either distribution.

. If one of the samples of text is flat, say  $q_i \in \mathscr{U}_{\mathcal{C}}$ , then

*ξ* = / (26)

For

£

$$p = c \sum_{i} p_i q_i = \% \sum_{i} p_i = 1 \times 1 = 1$$

If both samples are very rough, then  $\xi$  fluctuates widely as the q's are permuted. The question arises, how wide is the distribution of  $\xi$  ?

(25)

. .

A standard measure is the "variance"

$$\sigma^{2} = E(\xi^{2}) - [E(\xi)]^{2}$$
(27)

We evaluate 
$$E(\xi^2) = \frac{1}{2} (c \sum_{i}^{c} \rho_i q_i)^2$$
 (28)

where  $\sum_{T}^{r}$  is understood to mean the sum over all possible permutations of the q's. Then  $E(\xi^{2}) = \frac{1}{2}! \sum_{T} c^{2} \sum_{i=1}^{c} \sum_{j=1}^{c} \rho_{i} q_{i} \rho_{j} q_{j}$ (29)  $= \frac{1}{2}! \sum_{i=1}^{c} \sum_{j=1}^{c} \rho_{i} p_{j} \sum_{T} q_{i} q_{j}.$ The term  $\sum_{T} q_{i} q_{j}$  can be evaluated in each of the cases  $i \neq j$  and i = jWe will use  $P = c \sum_{i=1}^{c} \rho_{i}^{2}$  and  $Q = c \sum_{i=1}^{c} q_{i}^{2}$ 

as the I.C.'s of the two samples.  
For 
$$i \neq j \sum_{T} q_i q_j = (c-2)' \sum_{i \neq j}^{C} q_i q_j$$
  

$$= (c-2)! \sum_{l=1}^{C} q_l \sum_{i \neq j}^{C} q_j = (c-2)! \sum_{l=1}^{C} q_l (l-q_l)$$

$$= (c-2)! \sum_{l=1}^{C} (q_l - q_l^2) = (c-2)! (l-9c).$$
(30)

For 
$$i = j \sum_{\tau} q_i^2 = (c - l) l \sum_{i=l}^{c} q_i^2 = (c - l) l q_c^{\prime}$$
 (31)

Thus (29) becomes

$$E(\xi^{2}) = \zeta_{C}^{2} \sum_{i=1}^{c} \sum_{j \neq i}^{c} p_{i} p_{j} (c-2)! \frac{c-q}{c} + \zeta_{C}^{2} \sum_{i=1}^{c} p_{i}^{2} (c-1)! q_{C}^{2}$$

$$= \zeta_{C}^{2} (c-2)! \frac{c-q}{c} \sum_{i=1}^{c} \sum_{j \neq i}^{c} p_{i} p_{j} + \zeta_{C}^{2} (c-1)! q_{C}^{2} \sum_{i=1}^{c} p_{i}^{2}$$

$$= \zeta_{C}^{2} (c-2)! \frac{c-q}{c} \frac{c-p}{c} + \zeta_{C}^{2} (c-1)! q_{C}^{2} \frac{c-p}{c} = \frac{(c-q)}{c} \frac{(c-q)}{c} \frac{(c-p)}{c} + \frac{q-p}{c}$$

$$= \frac{c^{2}-cp-cq+pq}{c(c-1)} + \frac{cq}{c} \frac{c-p}{c(c-1)} + \frac{cq}{c(c-1)} + \frac{q-p}{c}$$

$$= \frac{c^{2}-cp-cq+pq}{c(c-1)} = \frac{c-p-q+pq}{c(c-1)}$$

$$(32)$$

Then since  $E(\xi) = I$  , we have

$$E(\xi^{2}) - [E(\xi)]^{2} = \frac{C - P - Q + PQ}{C - I} - I$$

$$= \frac{C - P - Q + PQ - C + I}{C - I} = \frac{PQ - P - Q + I}{C - I}$$

$$\sigma^{-2} = \frac{(P-I)(Q-I)}{C-I}$$
(33)

15

This says that the square of a standard deviation  $\sigma^-$  is the product of the bulges over c-1.

$$S = \frac{\xi - 1}{\sqrt{\frac{(P-1)(Q-1)}{C-1}}} = (\xi - 1) \sqrt{\frac{C-1}{(P-1)(Q-1)}}$$

a measure of the significance which can be judged from table I.

The function  $\xi$  can be used as a measure of correlation between two distributions.

#### XI <u>THE COINCIDENCE TEST USED TO ALIGN</u> SECONDARY ALPHABETS INTO A PRIMARY ALPHABET

.

We give here a special application of the I.C. statistic. An actual problem from the elementary course is used, problem 4 of assignment 6. Special frequency distribution tables were made of the sample lined up into 26 columns, i.e., lines of 26 letters each.

This problem happens to be enciphered by means of a Vigenere table, the columns being used in successions. Consequently if the cipher text is lined up 26 wide each column is enciphered by a monoalphabetic substitution. Each alphabet is a slide on that in the next column. If we knew the plain and cipher sequences the text could be decrypted. The problem is to recover these sequences. Since the sequence is not alphabetical, adjacent frequency counts as given in table C appear unrelated. But if we look at the rows they must be related by being slides on each other. If we can establish these slides will have the cipher sequence.

<u>Table "C"</u>	<b>Frequency</b>	Table	of	the	Cryptogram

Column - Cipher -	1	26	25	24	23	22	21	2Ø	19	18	17	16	15	14	13	12	11	1Ø	9	8	7	6	5	4	3	2	Total <sub>.</sub> Letters
A	1	3		3				2		5		3	2							1		2		1	2		25
Β.		1	3	1		4	4	1					2	5		5		2	3	1			1		1	2	36
C	<u> </u>	2		1	1							1		3	2			1	1			3			1	2	
D		1			1		1			2	1	1		1			7				4		3		1		<u>18</u> 22
E			1	2	2				2		3			1		1	1	1				1	4	5	1		25
F				1		<b></b>		1	1	2						1	1		2		1		1	2		1	14
G	1	1			1	5		1	1	2	1		1		2	_				1	1	1		2	3	2	26
H	3		2		1	1						2		1		2	1		1		1	1	1			1	<u>26</u> 17
Ī		1			2			1		1	1			1	1		2	1				3				4	18
J			3	1	3		3					1			2		1	4			1		•	1	2		22
· K	4	1			2			1	3			1	3	1	2	2		2			_	3		· · ·	2		27
L			1				4				3	2				3			1	-	3			1	1	1	20
M	1		1	3		1	1	1	2		3				5			1		1				1		1	22
N		1				t		3	3				4	3	1				1	2		2		2	2		22 24 16
0	3		2	1				1	1				2			2	1	1		-	1					1	16
P	-			3	3				1						1					4	4				1		17
Q •				1			1	<u> </u>			5	1			_		· ·	1	1	2			2	<u> </u>		2	16
R	1	1	1	2		1		2	4	1	1		2		1			3	1		1		1			2	24
S	2		3		1	2				4			2			1		-		2		1	3				24 21 26 19
T		2	1		1	<u> </u>						4	3	1			3			5			1	3	_	2.	26
Ū	1			1	1	4	1	3	1		2			1				2				2					19
v v	4	2	-		1			Ĩ	1			1		1	1	1		3	2		、	3	2		3	1	27
Ŵ	<u>↓</u>	4	1	1	<u>├─</u>	1	2	<u> </u>	-		1			_=	-	3	3		1	3	2			1	-		<u>13</u>
X	1	<u> </u>			1	<u> </u>	2	-		1		1		1			1		4	-	1		1	2	1		17
TY T	<u>† –</u>	1			+−	1	12	3	1	1	<u> </u>	2		-	1		<b></b>		3		1		$\hat{2}$	-	1	┝──╊	19
Z	1		2			Î	†- <u>-</u> -	<u> </u>	┢╼┈	2		1		1	2		1		1	<b>—</b>	1		-	1	<b></b>		15

,

1

.

By use of this special table, table "C", we can build up the cipher component used by matching these frequency distributions. In selecting distributions to match we want to obtain:

- (a) A maximum total number of letters involved (as the distribution will then be more reliable).
- (b) A distribution with a normal count (i.e., similar to normal alphabet).
- (c) A distribution without any one letter of abnormal frequency (as this gives too much weight to one letter).

When frequency distributions of two letters are properly matched, high should pair with high, low with low, blank with blank, etc. The mathematical value of each .relative position is found as follows:

(a) With the two frequency distributions in question written on paper strips and slid one against the other, for any one position we cross multiply the frequencies in alignment, and then add the products of all these multiplications. This is the total number of parts or coincidences involved (see table "D").

(b) Cross-multiply the total count of the distribution of the first letter by the count of the second letter. This is the total number of possible pairs of letters. Chance would produce one coincidence in twenty-six. Therefore, divide this product by twenty-six, which gives the number of expected chance coincidences.

(c) Divide the number of the actual coincidences by the expected number of chance coincidences (that is, divide (1) by (2)). The resulting number is the <u>Index of</u> <u>Coincidence</u>.

To prove correct alignment:

(a) The index-for the given relative position of two distributions must be higher than for all other positions, with no close second.

(b) The index should be 1.50 or higher (preferably 1.75 or higher).

(c) There must be only one acceptable alignment.

Indeterminate results will be encountered in some cases, particularly with insufficient text.

Table "E" gives the total councidences at the various positions of one strip slid against another.

From table "C", it is seen that certain distributions have the following properties (referring to our three desired properties):

B (cipher) has a total of 36 letters, with 14 different cells involved. Its highest frequency is 5. Good. Approaches normality. Maximum text.

V (cipher) has a total of 27 letters with 15 different cells involved. Its highest frequency is 4. Not good -- too flat.

K (cipher) has a total of 27 letters with 13 different cells involved. Its highest frequency 15 4. Not good - too flat.

G (cipher) has a total of 25 letters with 15 different cells involved. Its highest frequency is 5. Not good -- too flat.

D (cipher) has a total of 22 letters with only  $1\emptyset$  different cells involved, but its highest frequency is 7. Not good -- too peaked.

T (cipher) has a total of 26 letters with 11 different cells involved. Its highest frequency is 5. Good. Approaches normality.

A (cipher) has a total of 25 letters with 11 different cells involved. Its highest frequency is 5. Good. Approaches normality.

N (cipher) has a total of 24 letters with 11 different cells involved. Its highest frequency is 4. Good. Approaches normality.

B, T, A and N are the best choices. Match T, A and N against B, then match A and N against T, finally match N against A. One of these combinations should give a positive index of coincidence, and thus serve as a starting point.

.

æ

Table	''D''	Sliding Strip	S

.

•

Plain Cipher	:::		2	3	4	5	6	7	8	9 :	1Ø	11	12	13	14	15	16	17	18	19 ``	2Ø	21	22	23	24	25	26:	"	B''	Mast	ter	Di	stri	but	lon
Freq.	:		1	3	1		4	4	1					2	5		5		2	3	1			1		_1	2:		13	1	4	4	1	2	5
			т	2	1		1							4	3	1			3			5			1	3		2: :					2 on ribu		
•	: 1 : 4		3		3			,	2		5		3	2							1		2		1	2	:								
	: 1	4	1						3	3			1	4	3	1		• ••		1	2		2		2	2	:								
Plain Cipher	: :	1 r	2	3	4	5 N	6	7	8	9	1ø	11	12	13	14	15	16	17	18	19	2Ø	21	22	23	24	25				N ibut			aste	r	
Freq.	:		4	3		1	1						7	6	1			7	3	1	5			2	5		: 4:	4	3		1 1	<u>L</u>	7	6	1
												: : B	1	3	1		4	4	1					2	5		5	2	3	1	1	1	2:		
																	: :A	1	3		5				2		5	3	2				1 2	1 :	

ø

.

•

.

.

.

.

.

۰

۰.

~

.

.

`

1

# Table "E" -- Table of Total Coincidences

\$

- 20

$\mathbf{B} = \mathbf{M}$	las	ter	· I	)is <sup>.</sup>	tri	but	ion		-																	•		
Plain		1	2	3	4	ŧ	56	7	8	9	1Ø	11	12	13	14	15	16	17	18	19	2Ø	21	22	23	24	25	26	Chance Coinci
Ciphe	r-	B																										dences
T							)																					36
A		_																										-
N	_	Х -					- 53	49												43	39	63						33
		NOT	'E :		The	nı	ımbe	rs	43.	49	. e	tc.	. r	enro	esei	nt ·	the	SU	cce	ssi	ve							
							L co																					
							enci																•					
						•							•															
T = M	las	ter	·I	)is	tri	but	ion																					
		·	~	~					~	~			10	- 0			10		10	10	• •	01	~~	~~	0.4	05	00	
<u>Plain</u>	•		2	3	4	-	5 6	7	8	9	10	TT	12	13	14	19	10	17	19	19	20	21	22	23	24	25	20	
<u>Ciphe</u>				20					26	20		12				21		51		43					<u>.</u>			25
							y																				34	24
N	_	Λ -			20	,	<i>y</i>	20	20				JI	,			51			51						02	01	27
A = M	las	teı	: 1	)is	tri	bu	tion																					
	-		_	_					_	-															<b>.</b> .			
Plain			2	3	4	;	56	7	8	9	1Ø	1 <b>T</b>	12	13	14	15	16	17	18	19	2Ø	21	22	23	24	25	26	
Ciphe	<u>r</u> -	A												~ ~					~ ~					~ =	•••		~ =	
N	-	X -	-	26									28	26		43		34	27					27	29		27	23
<b>T</b> - 1			1	r _	E	_	Mo	-+-	- D				~~															
1 - 1	•		· 1		J	=,	n na J	sle	r D.	ISC.	CT 0	uti	on															
Plain	_	1	2	3	4		56	7	8	q	10	11	12	13	14	15	16	17	18	19	2Ø	21	22	23	24	25	26	
Ciphe	-	_	-	Ŭ	-	1		•	0	. 7	10			10		10	<i>,</i> • •	<b>±</b> 1	10	20	2,0	~.		20	- •			
							ζ	10	ø –			<b>`12</b>	3 -	8Ø	78	74						82		75	73	89	96	69
																												- 48
												~*																

REF ID:A66796

# Table "F" -- Table of Coincidences

# Master Distribution T-1 N-5 B-11 A-17

Plain	1	2	3	4	5	6	7	8	9	1Ø	11	12	13	14	15	16	<b>µ7</b>	18	19	2Ø	21	22	23	24	25	26	Chance
													~							<u> </u>					_		Coin-
Cipher	Т	V	Ι	D	N	Y	ប	F	P	M	B	K	H	Q	Z	Е	A	L	G	X	R	J	W	S	С	0	ciden-
-	Т				N						В						A	ł		1						{ }	<u>ces</u>
ĸ	X	149	-	-	X	-	-	-	-	-	X	175	-	-	-	-	X	- 1		-	-	-	-		-	156	120
v		181	-	_	X	-	-	-	-	-	X	X	-	-	-	-	X	-	-	-	-	143	-	-	-	158	120
Ď.	X	X	_	193	X	-	_	165		_	X	X	-	-	-	-	X	-	-	-	-	143	_	~	-	- 1	94
E	X	X	135		X	-	-		-	-	X	X	-	-	-	19Ø	X	- 1	-	-	148	-	-	-	-		102
W	X	X	_	-	X	-	-			129	X	X	-	-	-	X	X	-	-	-	-	~	200	-	-	1-1	102
G	X	X	-	-	X	-	-		-	-	X	X	-	-	143	X	X	- 1	193		131	-	X	-	-	1-1	107
R	X	X		-	X	-	-	1	133	-	X	X	-	-	-	X	X	- 1	X	-	158	-	X	-	-	- 1	102
M	X	X	133	-	X	-	-	-	-	155	X	X	-	144	-	X	X	-	X	-	X	-	X	-	1	- 1	94
J	X	X	-	-	X	-	-	-	-	-	X	X	-	-	-	X	X	-	X	-	X	1 7Ø	X		-	- 1	9Ø
L	X	X		-	X	144	-	-	1	-	Х	X	-	-	-	X	X	154	X	-	X	X	X	-	~	1-1	9Ø
S	X	X		-	X	-	-	110	1	-	X	X	-	-	109	X	X	-	X	- 1	. X	X	X	159	-		<u>9ø</u>
Ċ	X	X	117	-	X	-	_	-	1	-	X	X	-	-	-	X	X	-	X	-	X	X	X	X	172	- 1	81
U	X	X	-	-	X	-	119	9Ø	-	_	X	X	-	- 1	- 1	X	X	- 1	X	-	X	X	X	X	X	1-1	<u>81</u> 81
Ŷ	X	X	-	-	X	153	-	93	-	96	X	X	-	-	- 1	X	X	- 1	X	-	X	X	X	X	X	t <del>x t</del>	81
Ī	X	X	121	-	X	X	X	94	-		X	X	109	- 1	[	X	X	X	X	- 1	- 1	X	X	X	X	X	7,7
H	X	X	X	-	X	X	X	-	-	-	X	X	119	-	86	X	X	X	X	1	X	X	X	X	X	9Ø	73
P .	X	X	X	-	X	X	X	-	142	-	Х	X	X	-	-	X	X	X	X	116	X	X	X	X	X	1 - 1	73
x	X	X	X	-	X	X	X	11ø	X	-	X	X	X	-	-	X	Ťx	X	X	139	X	X	Χ.	X	X	1-1	73
0	X	X	X	- 1	X	X	X	-	X	-	X	-	X	X	-	-	- 1	X	X	X	X	X	X	X	X	117	68
Q	X	X	X	-	X	X	X	-	X	81	X	X	X	108	-	X	X	X	X	X	X	X	X	X	X	X	68
Z	X	X	X	-	X	X	X	-	X	X	Х	X	X	X	113	X	X	X	X	X	X	X	X	X	X	X	
F	X	X	X	72	X	X	X	87	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	<u>64</u> 6Ø

### BUILDING UP THE CIPHER COMPONENT

By utilizing the principles described in the previous sections, we can build up the cipher component. Take B (cipher) as the master distribution, since it is accepttable and contains the highest count. Copy the frequencies of B (from "C") at the bottom of a strip of paper, and repeat this sequence to the right. -Over the first sequence write the numbers 1 to 26 as shown in table "D". Under No. 1 write the letters"B". The numbers represent the various unknown letters of the cipher component. Make similar master distribution strips for T and A. Next, copy the distribution of T (cipher) (from table "D") at the top of a strip of paper. Only one sequence is required for this strip, and the numbers are omitted. Indicate the space corresponding to column No. 1 (table "D") by the letter T (see table "E"). In a like manner make strips for A and N.

<u>Note</u>: The letter on each strip is an indicator to mark column No. 1 for that letter. When strips are properly aligned, the indicators show the relative positions of these letters in the cipher component. The student is advised to prepare strips for himself and follow these processes.

First, match T against B. As no two letters can occupy the same position in the cipher component, begin by setting the T indicator at No. 2 on the B master alphabet. Note the coincidences. Next, slide T to No. 3, and note the coincidences. Continue this process to No. 26, and record the successive coincidences in tabular form (see table "F"). In many cases lack of good coincidences will be obvious by inspection tion and the count need not be made. In this way we discover that B and T give high indices of coincidence in two different alignments (indices computed in accordance with rule in page 23).

Index of B (1) - T (7)  $\frac{64}{36}$  = 1.77 (good) Index of B (1) - T (11)  $\frac{69}{36}$  = 1.67 (good)  $\frac{69}{36}$ 

All other alignments give such low indices that they can be at once eliminated. The above two indices, however, are both high enough to be significant, and as the second is so close to the first, it cannot be disregarded.

• There can be only one acceptable point of coincidence; therefore, it is necessary to match A against B, and N against B, to see if more conclusive results can be attained.

Index of B (1) - A (5)  $\frac{66}{36} = 1.89$  (excellent) 36Index of B (1) - A (7)  $\frac{59}{35} = 1.69$  (good) (other alignments are eliminated)

Index of B (1) - N (21)  $\frac{63}{33} = 1.91$  (excellent)'

Index.of B (1) - N (6)  $\frac{53}{33} = 1.6\emptyset$  (fair)

(other alignments are eliminated)

Since there is no outstanding coincidence with "B" as the "master-alphabet", try "T" and "A" as the "master-alphabet":

Index of T (1) - A (17)  $\frac{51}{25}$  = 2.04 (excellent) Index of T (1) - A (19)  $\frac{43}{25}$  = 1.72 (good) Index of T (1) - A (11)  $\frac{42}{25}$  = 1.68 (good) Index of T (1) - N (5)  $\frac{50}{24}$  = 2.08 (excellent) index of T (1) - N (12)  $\frac{34}{24}$  = 1.42 (poor) (16)  $\frac{24}{25}$  = 1.87 (good) Index of A (1) - N (17)  $\frac{34}{23}$  = 1.48 (poor)

٤

The most certain combination is  $T_{-}(1) - N$  (5), and there is no doubt as to its correctness. This located "N" relative to "T" in the cipher component, and allows us to consolidate their frequencies.

For a new master distribution add the frequencies of T (at space No. 1) to those of N (at space No. 5) (see table "E"). Match "B" and "A" against this new master distribution:

Index of T (1) - N (5) - B (11)  $\frac{123}{69} = 1.81$  (good) Index of T (1) - N (5) - B (7)  $\frac{100}{69} = 1.45$  (poor) Index of T (1) - N (5) - B (26)  $\frac{76}{69} = 1.40$  (poor) Index of T (1) - N (5) - A (17)  $\frac{94}{48} = 1.96$  (excellent) Index of T (1) - N (5) - A (15)  $\frac{65}{48} = 1.35$  (very poor) Index of T (1) - N (5) - A (19)  $\frac{65}{48} = 1.35$  (very poor)

"B" and "A" can now be consolidated with "T" and "N" for the final "master distribution" as follows:

													•
:Plain :Cipher	- 1 - T	2	3	4	5 N	6	7	8	9	10	11	12	13
:T &	_	4	3		1	1						7	6
: В		2	3	1			1		1	2		1	3
: A		3	2		_	<b>、</b>		·		1		2	
Plain	- 1	4 15	16	17	18	19	2Ø	21	22	23	24	25	26:
Cipher	-			-	~	1	-			2	e		
Т &			,		3	1	5			·2	5		4:
B			4	4	1		~			·Z	5		5:
A		1 2		1	3		3				2	•	5:

This locates "B" and "A" relative to "T" and "N" in the cipher component, in addition to giving the combined frequencies of all four letters.

:Plain :Comp.	T	1	2	3	4	.5	6	7	8	9	1ø	11	12	13
:Cipher	-	Т				N						В		
Comp. Comb. Freq.	-		9	8	1	1	1	1		1	3		1ø	9
Plain Comp.	-	14	15	16	17	18	19	2Ø	21	22	23	24	25	26
Cipher Comp.	-				A									:
Comb. Freq.	-	3	2	4	12	7	1	8		4	12			14

Analyze the preceding steps. T gave two possible alignments with B, and, as we now see, the incorrect position gave the higher index. N also gave two possible alignments with B. (B was at fault due to its erratic letter distribution). However, when T and N are combined, giving twice as many letters in the master distribution, B fitted in with only one possible alignment. Adding B and A gives twice as many letters in the master distribution and this should make future results even more positive. The master distribution (of 1000 letters or more) should approximate a normal frequency distribution and will give a standard to which all the other distributions can be referred. Hereafter, variations in the highest index of coincidence will be due entirely to letter distribution of the various distributions themselves.

For example:

If the highest index, is 1.7 -- letter distribution is normal.

FINAL MASTER DISTRIBUTION

If the highest index is  $2.\emptyset$  -- high frequency letters predominate.

If the highest index is 1.4 -- the intermediate and low frequency letters predominate.

Therefore, when matching the remaining letters, we can accept the <u>highest</u> index of coincidence as establishing coincidence, unless the second highest is practically the same.

Continue the matching process and the reconstruction of the cipher component, noting that T, N, B and A are already located and thus may be deleted at once from further test. Begin with the letters of the highest frequency, as they should give the most positive results. When a letter is placed, delete this location from further test, as two letters connot occupy the same space in the cipher component.

Letters are added to the cipher component in the following order (see table "F"):

"Master alphabet" T (1) - N (5) - B (11) - A (17)

- K (12) 1.46 (poor but acceptable)
- V (2) 1.51 (poor but acceptable)
- D (4) 2.05 (excellent)
- D (8) 1.76 (good)(D not certain)
- E (16) 1.86 (excellent)
- W (23) 1.96 (excellent)
- G (19) 1.80 (good)
- <u>Note</u>: With this many values a key-word (if any) sequence could be completed by inspection. In this case, the partially reconstructed cipher component gives no suggestion of a key-word sequence.

R	(21)	1.55	(fair bu	it acceptable)
			•	•

M (1Ø) 1.65 (good)

- M (14) 1.53 (M not certain)
- J (22) 1.89 (excellent)
- L (18) 1.71 (good)
- L (6) 1.6 $\emptyset$  (fair)(L not certain)
- S (24) 1.77 (good)
- C (25) 2.12 (excellent)
- U (7) 1.47 (poor but acceptable)
- Y (6) 1.89 (excellent)
- <u>Note</u>: This throws out L (6), but leaves L (18) as correct.
- I (3) (fair)
- I (13) 1.42 (poor)(I not certain)
- H (13) 4.63 (good)

Note: This throws out I (13) and leaves I (3) as correct.

- P (9) 1.95 (excellent)
- X (2Ø) 1.9Ø (excellent)
- 0 (26) 1.72 (good)
- Q (14) 1.59 (fair and acceptable)
- Note: This throws out M (14) and leaves M (10) as correct.
- Z (15) 1.77 (good)
- F (8) 1.45 (poor but acceptable)
- F (4) 1.20 (very poor)
- F (8) is correct and D (4) is correct

The cipher component has now been completely recovered.

Note: The process described above has actually built up the complete squaredcipher-table of a modified Vigenere table (it remains only to recover the plain component to complete the Vigenere table). We have written down the cipher component rather than the complete-squared-table merely to save time and effort.

### XII THE ROUGHNESS OF MIXED TEXTS

What happens when two different distributions are mixed? As a simple case, let us suppose we mix some text of I.C.  $\gamma$  with flat text in the proportions R:(1 - R). Then the ith letter has probability,

Then the left left has probability,  

$$p_{i}^{R} + (1 - R) \frac{1}{c}$$
and the I.C. is  $c \sum_{i}^{C} (\rho_{i} R + (1 - R) \frac{i}{c})^{2}$ 

$$= C \sum_{i}^{C} \rho_{i}^{2} R^{2} + 2c \sum_{i}^{C} \rho_{i} R (1 - R) \frac{i}{c} + c \sum_{i}^{C} (1 - R)^{2} \frac{i}{c^{2}}$$

$$=R^{2}\gamma + 2R(I-R) + (I-R)^{2}$$

 $= R^{2} \gamma - R^{2} + I = I + (\gamma - I) R^{2}.$ 

- 25 -

(34)

That is, the rough text contributes its bulge in the proportion  $\mathbb{R}^2$ .

Now as a more complicated case suppose that we mix two rough texts in the proportions R: (1 - R). Suppose further that the distributions are  $P_1$ ,  $P_2$ , ---,  $P_c$  and  $q_1$ ,  $q_2$ , ---,  $q_c$  and  $\sum_{n=1}^{n} p_n q_n q_n$ 

$$c \sum p_i^{z} = P, \qquad c \sum q_i^{z} = Q$$

Then the mixture will have as probability of the ith letter  $p_i R + q_i (l-R)$ and  $\gamma = c \sum (p_L R + q_i (l-R))^2$ 

$$= c \sum p_i^{2} R^{2} + 2c \sum p_i q_i R (1-R) + c \sum q_i^{2} (1-R)^{2}$$

$$= PR^{2} + 2R(I-R) c \sum \bar{p_{i}} q_{i} + Q(I-R)^{2}$$

The expression  $C \sum_{p_i} q_i = \mathcal{E}$  we have examined before, and seen to be a measure of the correlation of the distributions. Since the expected value of  $\mathcal{E}$  is 1, the expected value of  $\mathcal{P}$  will be  $E(\gamma) = PR^2 + 2R(1-R) + Q(1-R)^2$ 

$$= R^{2} + (P - I)R^{2} + 2R - 2R^{2} + (I - R)^{2} + (Q - I)(I - R)^{2}$$
  
= I + (P - I)R<sup>2</sup> + (Q - I)(I - R)<sup>2</sup>.

i=1

Here again each contributes its bulge in proportion to the square of ite weight.

Now that we have seen how the argument goes, we can generalize this result to a mixture of K samples with distributions  $P_{i1}$ ,  $P_{i2}$ , ---,  $P_{ic}$ , and  $\gamma_i = c$   $\sum_{j=1}^{c} p_{ij}^2$  the I.C. of the ith sample, which we will suppose present in the proportion  $R_1^i$ ,  $\sum_{l=1}^{K} R_l = l$ . The mixture will have the jth letter present in the proportion

$$Y = c \sum_{j=1}^{c} P_{j}^{2} = c \sum_{j=1}^{c} \left( \sum_{l=1}^{k} R_{l} - P_{lj} \right)^{2}$$
  
=  $c \sum_{j=1}^{c} \sum_{l=1}^{k} \sum_{L=1}^{k} R_{l} - P_{lj} - R_{L} - P_{Lj}$   
=  $\sum_{l=1}^{k} \sum_{L=1}^{k} R_{l} - R_{L} - C - \sum_{l=1}^{c} P_{lj} - P_{L}$ 

 $P_{j} = \sum_{l=1}^{n} R_{l} P_{lj}$ 

(36)

(35)

The internal sum is the measure of the correlation of the ith and  $\ell$  th samples, which we will designate by

Notice that 
$$\xi_{i\,i}^{\prime} = \gamma_{i}^{\prime}$$
 is the I.C. of the ith sample.  
Now  $\gamma$  becomes  $\gamma = \sum_{i=1}^{\kappa} \sum_{L=1}^{\kappa} R_{i} R_{L} \xi_{i\,L}^{\prime}$   
 $= \sum_{i=1}^{\kappa} \sum_{L\neq i}^{\kappa} R_{i} R_{L} \xi_{i\,L} + \sum_{i=1}^{\kappa} R_{i}^{\prime \prime} \gamma_{i}^{\prime}$ 
 $= \sum_{i=1}^{\kappa} \sum_{L\neq i}^{\kappa} R_{i} R_{L} \xi_{i\,L} + \sum_{i=1}^{\kappa} \left[ R_{i}^{\prime \prime} + (\gamma_{i}^{\prime} - I) R_{i}^{\prime \prime} \right]$ 
(37)

we have

(38)

Since  $E(\xi_{i_{L}}) = 1$ 

 $E(\gamma) = \sum_{l=1}^{\kappa} \sum_{l=1}^{\kappa} R_{l} R_{l} + \sum_{i=1}^{\kappa} (\gamma_{l} - 1) R_{l}^{2}$   $= \sum_{l=1}^{\kappa} R_{l} \sum_{l=1}^{\kappa} R_{l} + \sum_{l=1}^{\kappa} (\gamma_{l} - 1) R_{l}^{2}$   $E(\gamma) = 1 + \sum_{l=1}^{\kappa} (\gamma_{l} - 1) R_{l}^{2}$ 

Again each sample contributes its bulge according to the square of its presence.