# Use of Bayes Factors With a
# Composite Hypothesis

BY FRANCIS T. LEAHY

*Top Secret Dinar*

There is a message in a sealed envelope. I am told that it is either a message that was enciphered by a columnar transposition of plain text, or that it was enciphered on a machine with an irregularly stepping rotor maze. For the moment, we presuppose that the language of the deciphered message will be English.

Thus, I have two a priori hypotheses, I and II.

Next, the envelope is opened, and a test is to be made to determine which hypothesis is correct. For this purpose, a "Bayes Factor" is to be computed.

For every character in the message, a "weight" (i.e., a particular positive or negative number) is to be substituted, and the sum of all these weights then obtained. (If a frequency count of the message exists, a short cut would be to multiply the frequency of "A" by the weight of "A", the frequency of "B" by the weight of "B", etc., throughout the entire alphabet, and then to sum.)

The weights had previously been prepared with the thought in mind that they were going to be *added* together. Hence, the weight for each letter is the logarithm of the ratio of the probabilities of observing the letter in question under the first and under the second hypothesis.

The sum of the weights is therefore the logarithm of the Bayes Factor. When the antilog is obtained, it represents the *Odds* in favor of Hypothesis I rather than Hypothesis II.

In other words, a Bayes Factor is synonymous with Odds; strictly speaking, a posteriori odds. (*Note:* A "factor of a thousand-to-one" in favor of Hypothesis I is a phrase that is easily understood; but a "factor of one-thousandth-to-one" in favor of Hypothesis I had better, for the sake of clarity, be rephrased as a "factor of a thousand-to-one" in favor of Hypothesis II, even though this be an equivalent statement.)

1

It is important to remember that Bayes Factors always represent "Odds-in-favor" whereas race-track odds (with which the reader is presumably more familiar) are "Odds-Against." If the tote-board shows some horse's odds as 60–1, these may be interpreted as the odds in favor of his *losing* the race!

Now, suppose that the message just described was originally believed to be in one of three languages: English, French, or German. The identification of the language is not desired at the moment, merely the mode of encipherment; that is, by transposition or by maze. A Bayes Factor, of course, must be computed.

Weights should be prepared using ratios of the probabilities of all twenty-six (in this instance) characters under all six pairs of hypotheses:

|  | *Transposition* | | | *Rotor-Maze* | | |
|---|---|---|---|---|---|---|
| Hypothesis A | English | French | German | English | French | German |
| Hypothesis B | Flat | Flat | Flat | Flat | Flat | Flat |

"Flat," i.e., Flat random, the null hypothesis, assumes that all (26) categories are equiprobable; in this instance, each character has a probability of $\frac{1}{26}$. There will be 156 weights.

*Note:* A "reflecting" rotor-maze has cipher probabilities that differ from flat random; but if a "straight-through" maze had been presupposed, the cipher distribution of the characters would be flat random. This latter supposition would make all the weights equal to "0" (the log of 1) on the right half of the above figure.

Let us assume that base-10 logarithms are used in preparing the weights, although other bases are equally acceptable.

Now, using subscripts $(E, F, G)$ to show the language, and $(T, R)$ to show the mode, let $S_{ET}, S_{FT}, S_{GT}, S_{ER}, S_{FR}, S_{GR}$, represent the scores obtained after the weights are summed. These are the log Bayes Factors for each of the six sets of sub-hypotheses, A and B. Let the symbol O $(T/FR)$ be the odds in favor of transposition rather than flat random. O $(R/FR)$ represents the odds in favor of the rotor-maze rather than flat random. O $(T/R)$ are the odds in favor of transposition rather than rotor-maze.

$$\text{Then } O\ (T/FR) = \tfrac{1}{3}\ (10^{S_{ET}}) + \tfrac{1}{3}\ (10^{S_{FT}}) + \tfrac{1}{3}\ (10^{S_{GT}})$$

$$O\ (R/FR) = \tfrac{1}{3}\ (10^{S_{ER}}) + \tfrac{1}{3}\ (10^{S_{FR}}) + \tfrac{1}{3}\ (10^{S_{GR}})$$

$$\text{Finally, } O(T/R) = \frac{O\ (T/FR)}{O\ (R/FR)}\ .$$

The fractions ($\frac{1}{3}$) appearing before the brackets represent the a priori probabilities that the message was in English, French, or German respectively. Any other appropriate fractions that sum to 1, say $\frac{1}{4}$, $\frac{1}{2}$, $\frac{1}{4}$, representing the a priori probabilities of these three sub-hypotheses might have been employed.

The six sub-hypotheses have been reduced to two "composite" hypotheses.

Two rules must always be followed when dealing with composite hypotheses:
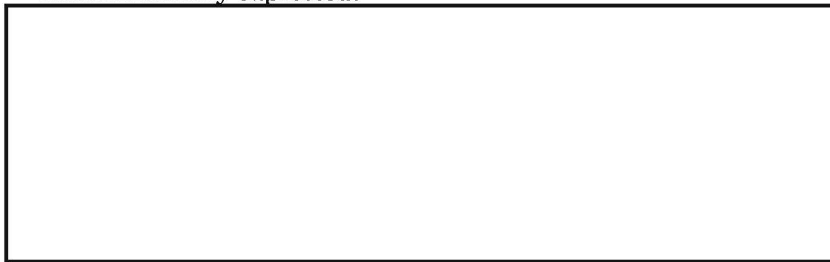
    1. The alternate sub-hypotheses must be the "null" or flat-random hypotheses.

    2. The probabilities of each of the several variations of the basic hypotheses must be cross-multiplied by the appropriate *odds*, never the *log-odds*.

At the very end, two distinct composite hypotheses, each of which has been compared against flat random and odds determined, can be compared with each other to determine their relative odds.

A practical example, much more complex than the one just illustrated, might run along the following lines:

EO 1.4.(c)
P.L. 86-36

Mathematically expressed:

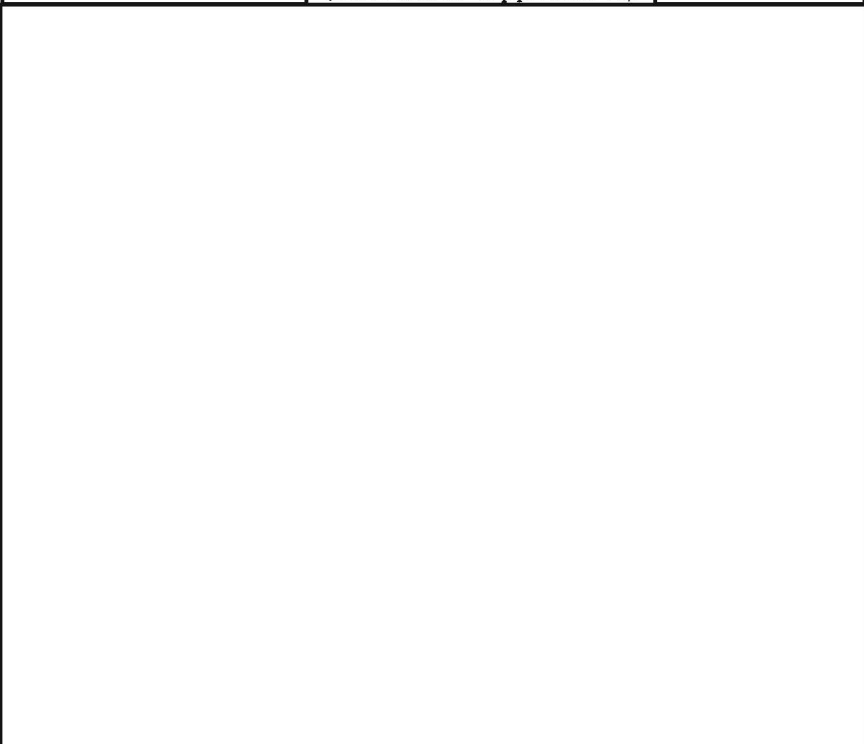The identical procedure is repeated for the second message.

It is obvious that *if* the odds that have just been computed had resulted in virtual certainties (say 100-to-1 in favor of active, 200-to-1 in favor of inactive, etc.), the superposition of the two sets of odds, at the right offset, would be "visible to the naked eye."
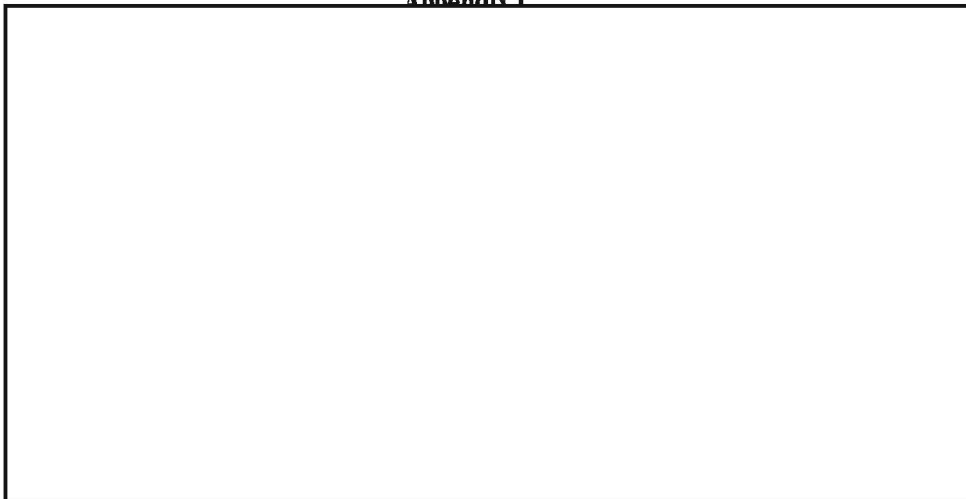
But, dealing with a much weaker case, which would arise whenever messages are short, it would still be possible to determine the *relative* probabilities of the UNCODED offsets. (Formula in Appendix I.)

EO 1.4.(c)
P.L. 86-36

FRANCIS T. LEAHY        ~~TOP SECRET DINAR~~

**Appendix I**

EO 1.4.(c)
P.L. 86-36

This is a good example of a Bayes Factor with a composite hypothesis.

EO 1.4.(c)
P.L. 86-36