

Frequency distribution activity

Goal: Turn the following cipher text into plaintext, using frequency distribution. It is known that the ciphertext and plaintext are both written using the English language.

yfzh zh mj kpktlzhk zj etkmizjq lzwfkth dhzjq vtkcdkjlq bzhytzedyzj. zy zh kmhzkhy yn dhk yfzh skyfnb zv gnd zjytklkwy m skhhmqk yfmy zh sntk yfmj njk fdjbtbb lfmtmlykth zj akjqyf, eklmdhk zj yfmy lmhk, yfk hmswak hzuk zh amtqk kjndqf hn yfmy yfk lfmtmlykt bzhytzedyzj zh sntk azikag yn lanhkag tkhkseak yfmy nv m ygwzlma skhhmqk xtzyykj zj kjqazhf. xk fmok smbk gndt ymhi kmhzkt zj yfzh kpmswak eg hkwmymyzjq yfk xntbh zj yfk skhhmqk.

The message is long enough to warrant attempting decryption using frequency distribution. The first step is to make a list of the frequency distribution of the letters in the ciphertext. It is as follows:

a - 11

b - 7

c - 1

d - 10

e - 6

f - 21

g - 7

h - 35

i - 3

j - 21

k - 52

l - 12

m - 29

n - 15

o - 1

p - 2

q - 10

r - 0

s - 10

t - 21

u - 1

v - 3

w - 6

x - 3

y - 35

z - 35

The following is an approximation of the distribution of letters in English, given a random writing sample of 1000 characters:

A - 73

B - 9

C - 30

D - 44

E - 130

F - 28

G - 16

H - 35

I - 74

J - 2

K - 3

L - 35

M - 25

N - 78

O - 74

P - 27

Q - 3

R - 77

S - 63

T - 93

U - 27

V - 13

W - 16

X - 5

Y - 19

Z - 1

In the ciphertext above, the letter k appears most frequently, with 52 instances, and the letters h, y and z are next, with 35 instances each. This implies that the cipher "k" probably translates to the plaintext "e," since "e" is the most commonly-found letter in English, and the cipher letters "h," "y" and "z" probably translate into three of the letters "a," "i," "n," "o," "r," "s" or "t," since these are the next-most-common letters in English.

Changing the ciphertext "k" to the plaintext "E," we have the following:

yfzh zh mj EpEtIzhE zj etEmizjq lzwfEth dhzjq vtEcdEjlg bzhytzedyzjn. zy zh EmhzeHy yn dhE yfzh sEyfnb zv gnd zjyEtIEwy m sEhmqE yfmy zh sntE yfmj nJE fdjbtEb lfmtmlyEth zj aEjqyf, eElmdhE zj yfmy ImhE, yfE hmswaE hzuE zh amtqE Ejndqf hn yfmy yfE lfmtmlyEt bzhytzedyzjn zh sntE aziEag yn lanhEag tEhEseaE yfmy nv m ygwzIma sEhmqE xtzyyEj zj Ejqazhf. xE fmoE smbE gndt ymhi EmhzeEt zj yfzh EpmswaE eg hEwmtmyzjq yfE xntbh zj yfE sEhmqE.

A study of short words (two or three letters) comes in handy here. We notice a few patterns, especially in reference to the most common letters seen in this ciphertext. Since “z” and “h” appear so frequently, and we notice there are five instances of the two-letter word “zh,” a good guess is that “zh” could be “IS,” “IN,” “AT,” “AN” or “OR.” Also, there are six times where “zj” appears, giving more strength to this argument. Let’s try the cipher “z” corresponding to the plaintext “I,” with the cipher “h” corresponding to the plaintext “S” and the cipher “j” corresponding to the plaintext “N.” Then we have:

yfIS IS mN EpEtIISE IN etEmiINq llwfEtS dSINq vtEcdENlg bISytlEdyInN. Iy IS EmSIESy yn dSE yfIS sEyfnb Iv gnd INyEtIEwy m sESSmqE yfmy IS sntE yfmN nNE fdNbtEb lfmtmlyEtS IN aENqyf, eElmdSE IN yfmy ImSE, yfE SmswaE SluE IS amtqE ENndqf Sn yfmy yfE lfmtmlyEt bISytlEdyInN IS sntE aliEag yn lanSEag tESEseaE yfmy nv m ygwllma sESSmqE xtIyyEN IN ENqalSf. xE fmoE smbE gndt ymSi EmSIEt IN yfIS EpmswaE eg SEwmtmyINq yfE xntbS IN yfE sESSmqE.

Looking at the first two words, “yfIS IS,” one might guess that this means “this is,” especially with the cipher “y” appearing 35 times. Guessing that provides us with:

THIS IS mN EpEtIISE IN etEmiINq llwHEtS dSINq vtEcdENlg bISTtledTInN. IT IS EmSIESt Tn dSE THIS sETHnb Iv gnd INTEtIEwT m sESSmqE THmT IS sntE THmN nNE HdNbtEb IHmtmITetS IN aENqTH, eElmdSE IN THmT ImSE, THE SmswaE SluE IS amtqE ENndqH Sn THmT THE IHmtmITet bISTtledTInN IS sntE aliEag Tn lanSEag tESEseaE THmT nv m Tgwllma sESSmqE xtITTEN IN ENqalSH. xE HmoE smbE gndt TmSi EmSIEt IN THIS EpmswaE eg SEwmtmTINq THE xntbS IN THE sESSmqE.

Looking at the second line, we find a one-letter word “m.” Since the plaintext “I” is already used, this must mean the cipher “m” corresponds to the plaintext “A.” This gives:

THIS IS AN EpEtIISE IN etEAiINq llwHEtS dSINq vtEcdENlg bISTtledTInN. IT IS EASIESt Tn dSE THIS sETHnb Iv gnd INTEtIEwT A sESSAQE THAT IS sntE THAN nNE HdNbtEb IHAtAITetS IN aENqTH, eEIAdSE IN THAT IASE, THE SAswaE SluE IS aAtqE ENndqH Sn THAT THE IHAtAITet bISTtledTInN IS sntE aliEag Tn lanSEag

tESEseaE THAT nv A TgwIIAa sESSAQE xtITTEN IN ENqaISH. xE HAoE sAbE gndt TASI EASIEt IN THIS EpAswaE eg SEwAtATINq THE xntbS IN THE sESSAQE.

In the second line, there is a two-letter word, "Tn," which implies that the cipher "n" is the plaintext "O." Then, in the fourth line, the word "ENndqH" implies "ENOUGH," so that the cipher "d" is the plaintext "U," and the cipher "q" is the plaintext "G." When we put in those three substitutions, we get:

THIS IS AN EpEtIISE IN etEAIING IiwHEtS USING vtEcUENIlg bISTtIeUTION. IT IS EASIESt TO USE THIS sETHOb Iv gOU INTetIewT A sESSAGE THAT IS sOtE THAN ONE HUNbtEb IHAtAITEtS IN aENGTH, eEIAUSE IN THAT IASE, THE SAswaE SluE IS aAtGE ENOUGH SO THAT THE IHAtAITEt bISTtIeUTION IS sOtE aliEag TO laOSEag tESEseaE THAT Ov A TgwIIAa sESSAQE xtITTEN IN ENGaISH. xE HAoE sAbE gOUt TASI EASIEt IN THIS EpAswaE eg SEwAtATING THE xOtbs IN THE sESSAGE.

It's getting much easier now, because we can see obvious words formed. For instance, in the third line, "HUNbtEb" implies "HUNDRED," in the fifth line, "ENGaISH" implies "ENGLISH," and in the last line, "sESSAGE" implies "MESSAGE." When we make those substitutions, we see:

THIS IS AN EpERIISE IN eREAIING IiwHERS USING vREcUENIlg DISTRIeUTION. IT IS EASIESt TO USE THIS METHOD Iv gOU INTERIEwT A MESSAGE THAT IS MORE THAN ONE HUNDRED IHARAItERS IN LENGTH, eEIAUSE IN THAT IASE, THE SAMwLE SluE IS LARGE ENOUGH SO THAT THE IHARAItER DISTRIeUTION IS MORE LIIElg TO ILOSElg RESEMeLE THAT Ov A TgwIIAl MESSAGE xRITTEN IN ENGLISH. xE HAoE MADE gOUR TASI EASIER IN THIS EpAMwLE eg SEwARATING THE xORDS IN THE MESSAGE.

Although we have deduced barely more than half the letters of the alphabet so far (14, to be exact), we have deciphered the vast majority of the letters in the ciphertext, and in fact, the rest is almost trivial. The cipher "l" obviously turns into the plaintext "C," and with that, things become clearer still, as shown here:

THIS IS AN EpERCISE IN eREAIING CIwHERS USING vREcUENcG DISTRIeUTION. IT IS EASIESt TO USE THIS METHOD Iv gOU INTERCEwT A MESSAGE THAT IS MORE THAN ONE HUNDRED CHARACTERS IN LENGTH, eECAUSE IN THAT CASE, THE SAMwLE SluE IS LARGE ENOUGH SO THAT THE CHARACTER DISTRIeUTION IS MORE LIIElg TO CLOSElg RESEMeLE THAT Ov A TgwICAL MESSAGE xRITTEN IN ENGLISH. xE HAoE MADE gOUR TASI EASIER IN THIS EpAMwLE eg SEwARATING THE xORDS IN THE MESSAGE.

Rather than go through the rest of the letters step-by-step, let's look at the message in its entirety:

THIS IS AN EXERCISE IN BREAKING CIPHERS USING FREQUENCY DISTRIBUTION. IT IS EASIEST TO USE THIS METHOD IF YOU INTERCEPT A MESSAGE THAT IS MORE THAN ONE HUNDRED CHARACTERS IN LENGTH, BECAUSE IN THAT CASE, THE SAMPLE SIZE IS LARGE ENOUGH SO THAT THE CHARACTER DISTRIBUTION IS MORE LIKELY TO CLOSELY RESEMBLE THAT OF A TYPICAL MESSAGE WRITTEN IN ENGLISH. WE HAVE MADE YOUR TASK EASIER IN THIS EXAMPLE BY SEPARATING THE WORDS IN THE MESSAGE.