

Models of Natural Language

In modeling English (or another natural language) we ignore *meaning* and even grammar. That is, we ignore **semantics**.

Instead, we treat English as being the result of some random process.

The worst possible model would be to suppose that each letter a-z in a message was chosen at random with **equal probabilities**. If this were so, cryptanalytic attacks would be nearly impossible. Luckily, this is not a good model.

Instead, the next-to-dumbest model observes that the **single-character** (and other) stats in English and other natural languages are *not flat*, but are *skewed* in a way that is useful.

In English, ignoring spaces (about 17% otherwise) the frequencies are roughly

e	occurs about	11 %	of the time
t	occurs about	9 %	of the time
o	occurs about	8 %	of the time
i	occurs about	8 %	of the time
a	occurs about	8 %	of the time
n	occurs about	7 %	of the time
s	occurs about	7 %	of the time
...			
v	occurs about	1 %	of the time
k	occurs about	1 %	of the time
x	occurs about	0.3 %	of the time
j	occurs about	0.23 %	of the time
q	occurs about	0.12 %	of the time
z	occurs about	0.09 %	of the time

Other natural languages using single-letter alphabets have similar skewings away from a *flat* distribution.

Digrams are *adjacent pairs* of characters. There are $26 \times 26 = 676$ different digrams. Every letter of the alphabet really does occur in English, but some digrams never occur.

In a sample of 500k of filtered email *with spaces removed*, 659 of the possible 676 digrams occur.

Top 44 digrams give more than 50% of the total.

The top 54 give 50%, the top 126 give 75%, the top 222 give 90%, and the top 359 give 98%.

<i>th</i> 3.18	<i>in</i> 2.59	<i>he</i> 2.17	<i>er</i> 1.95
<i>re</i> 1.85	<i>on</i> 1.63	<i>an</i> 1.59	<i>at</i> 1.54
<i>ou</i> 1.43	<i>or</i> 1.26	<i>es</i> 1.26	<i>ha</i> 1.24
<i>to</i> 1.22	<i>te</i> 1.21	<i>is</i> 1.18	<i>ti</i> 1.17
<i>it</i> 1.16	<i>en</i> 1.13	<i>nt</i> 1.09	<i>ng</i> 1.08
<i>al</i> 1.07	<i>se</i> 1.05	<i>st</i> 1.01	<i>nd</i> 0.98
<i>le</i> 0.91	<i>ar</i> 0.90	<i>me</i> 0.90	<i>hi</i> 0.86
<i>ve</i> 0.85	<i>of</i> 0.84	<i>ed</i> 0.78	<i>co</i> 0.74
<i>as</i> 0.73	<i>ll</i> 0.72	<i>ne</i> 0.70	<i>om</i> 0.70
<i>ri</i> 0.68	<i>ic</i> 0.67	<i>ro</i> 0.67	<i>ea</i> 0.66
<i>et</i> 0.64	<i>ur</i> 0.64	<i>io</i> 0.64	<i>ra</i> 0.62
<i>li</i> 0.62	<i>no</i> 0.62	<i>so</i> 0.62	<i>be</i> 0.61
<i>de</i> 0.59	<i>ma</i> 0.59	<i>si</i> 0.58	<i>ly</i> 0.54
<i>ut</i> 0.53	<i>ot</i> 0.53	<i>pr</i> 0.53	<i>fo</i> 0.53
<i>yo</i> 0.52	<i>il</i> 0.50	<i>ca</i> 0.50	<i>pe</i> 0.50
<i>ch</i> 0.49	<i>ho</i> 0.49	<i>ul</i> 0.47	<i>ce</i> 0.47
<i>ta</i> 0.45	<i>di</i> 0.45	<i>rs</i> 0.45	<i>el</i> 0.44
<i>ge</i> 0.44	<i>us</i> 0.44	<i>ec</i> 0.42	<i>ss</i> 0.42
<i>ac</i> 0.41	<i>ct</i> 0.41	<i>em</i> 0.41	<i>wh</i> 0.40
<i>oo</i> 0.40			

versus $1/656 \sim .15\%$

Trigrams are *adjacent triples* of characters. Of $26 \times 26 \times 26 = 17,576$ possible few occur often. The top 240 give 50% of all trigrams occurring: 1/70 of all trigrams account for 50% of occurrences.

When blanks are removed frequencies are spread out (as with digrams). The top 430 give 50%, the top 1160 give 75%, the top 2300 give 90%, and the top 4400 give 98%.

<i>the</i> 2.44	<i>ing</i> 1.26	<i>and</i> 0.82	<i>hat</i> 0.78
<i>tha</i> 0.77	<i>ion</i> 0.75	<i>you</i> 0.67	<i>ent</i> 0.66
<i>for</i> 0.63	<i>tio</i> 0.63	<i>thi</i> 0.60	<i>her</i> 0.51
<i>ati</i> 0.47	<i>our</i> 0.47	<i>ere</i> 0.45	<i>all</i> 0.43
<i>ter</i> 0.43	<i>ver</i> 0.40	<i>not</i> 0.40	<i>hin</i> 0.40
<i>ome</i> 0.36	<i>oul</i> 0.36	<i>uld</i> 0.36	<i>int</i> 0.34
<i>rea</i> 0.34	<i>pro</i> 0.34	<i>res</i> 0.33	<i>ate</i> 0.33
<i>hav</i> 0.30	<i>ave</i> 0.30	<i>ill</i> 0.30	<i>his</i> 0.30
<i>com</i> 0.30	<i>ons</i> 0.30	<i>are</i> 0.28	<i>ple</i> 0.28
<i>ers</i> 0.28	<i>con</i> 0.27	<i>ess</i> 0.27	<i>out</i> 0.27
<i>one</i> 0.26	<i>ith</i> 0.25	<i>som</i> 0.25	<i>ive</i> 0.25
<i>tin</i> 0.25	<i>nce</i> 0.24	<i>ble</i> 0.24	<i>ted</i> 0.24
<i>han</i> 0.23	<i>ine</i> 0.23	<i>per</i> 0.23	<i>ect</i> 0.23
<i>nte</i> 0.23	<i>wit</i> 0.22	<i>men</i> 0.22	<i>but</i> 0.22
<i>wou</i> 0.21	<i>ica</i> 0.21	<i>eve</i> 0.21	<i>cal</i> 0.21
<i>pre</i> 0.21	<i>cou</i> 0.21	<i>lin</i> 0.21	<i>est</i> 0.20
<i>eri</i> 0.20	<i>mor</i> 0.20	<i>ser</i> 0.20	<i>ore</i> 0.19
<i>any</i> 0.19	<i>abl</i> 0.19	<i>tic</i> 0.19	<i>urs</i> 0.19
<i>ant</i> 0.19	<i>sti</i> 0.18	<i>ear</i> 0.18	<i>hou</i> 0.18
<i>ies</i> 0.18			

versus $1/26^3 \sim 0.00569\%$

There are few **small words** in English, so if word breaks are preserved cryptanalysis is vastly easier.

Only 3 single-letter words, I and a

Only about 35 two-letter words
(versus $26^2 = 676$)

Only about 200 common three-letter words
(versus $26^3 = 17576$)

Common words are indeed very common. In 500k of my email, more than 5000 words appear, but the 9 most common words are 21% of all occurrences, the 20 most common give 30%, the 104 most common give 50%.

<i>the</i> 4.65	<i>to</i> 3.02	<i>of</i> 2.61
<i>i</i> 2.2	<i>a</i> 1.95	<i>and</i> 1.82
<i>is</i> 1.68	<i>that</i> 1.62	<i>in</i> 1.57
<i>it</i> 1.22	<i>for</i> 1.17	<i>you</i> 1.06
<i>be</i> 0.99	<i>not</i> 0.84	<i>on</i> 0.76
<i>have</i> 0.71	<i>this</i> 0.69	<i>as</i> 0.57
<i>at</i> 0.56	<i>would</i> 0.55	<i>are</i> 0.55
<i>but</i> 0.54	<i>if</i> 0.53	<i>my</i> 0.53
<i>with</i> 0.5	<i>your</i> 0.48	<i>so</i> 0.48
<i>or</i> 0.46	<i>some</i> 0.43	<i>will</i> 0.41
<i>do</i> 0.39	<i>about</i> 0.39	<i>me</i> 0.38
<i>from</i> 0.35	<i>by</i> 0.33	<i>no</i> 0.33
<i>more</i> 0.33	<i>what</i> 0.32	<i>an</i> 0.32
<i>there</i> 0.32	<i>one</i> 0.32	<i>all</i> 0.32
<i>was</i> 0.30	<i>we</i> 0.30	<i>just</i> 0.27
<i>which</i> 0.27	<i>can</i> 0.26	<i>very</i> 0.25
<i>series</i> 0.25	<i>am</i> 0.24	<i>things</i> 0.24
<i>people</i> 0.24	<i>get</i> 0.23	<i>hi</i> 0.23
<i>time</i> 0.22	<i>think</i> 0.22	<i>course</i> 0.22
<i>etc</i> 0.22	<i>also</i> 0.21	<i>any</i> 0.21

These low-level statistical biases can be put to use in cryptanalysis. On the other hand, they must be masked to have a secure cryptosystem.

A related point is that *non-random data can be **compressed***, which among its other effects *flattens* the stats of the code alphabet.

Just as it is good to remove spaces before encryption, it might be good to *compress* before encryption.

But, actually, good ciphers are so good that this does not really matter.

And, compressed data has highly structured, non-random **headers**, which partly defeat the purpose of flattening the stats.

Cryptograms

Cryptograms are also **monoalphabetic substitution ciphers**.

Substitution cipher means that each letter of the alphabet a, b, c, \dots, x, y, z is encrypted as some *other* letter by a secret chosen bijective function from the alphabet to itself.

Monoalphabetic means that the encryption of each letter is the same throughout the encryption of a message (rather than possibly varying over the message).

With spaces left in and with hints, they are in newspapers as puzzles. They are not secure.

Shift ciphers and affine ciphers are examples of monoalphabetic substitution ciphers in which the encryption function is particularly simple.

The typical non-machine setup
for a cryptogram is to have a key
which is a secret phrase, such as
My dog has fleas

Then remove duplicate letters, spaces, and
punctuation

mydoghasfle

Then put unused letters of the
alphabet on the end: append

bcijknpqrtuvwxyz

to get

mydoghasflebcijknpqrtuvwxyz

Then make the lookup table

abcdefghijklmnopqrstuvwxyz
mydoghasflebcijknpqrtuvwxyz

where the character in the top line is
encrypted as the character under it in the
lower line.