# Elementary Cryptanalysis

The most direct attack on a cryptosystem is an *exhaustive key search* attack. The key size therefore provides a lower bound on the security of a cryptosystem. As an example we compare the key sizes of several of the cryptosystems we have introduced so far. We assume that the alphabet for each is the 26 character alphabet.

Substitution ciphers:
    Simple substitution ciphers: 26!
    Affine substitution ciphers: $\varphi(26) \cdot 26 = 12 \cdot 26 = 312$
    Translation substitution ciphers: 26
Transposition ciphers:
    Transposition ciphers (of block length $m$): $m!$
Enigma :
    Rotor choices (3 of 5): 60
    Rotor positions: $26^3 = 17576$
    Plugboard settings: 105578918576
    Total combinations: 111339304373506560

The size of the keyspace is a naive measure, but provides an upper bound on the security of a cryptosystem. This measure ignores any structure, like character frequencies, which might remain intact following encryption.

# Classification of Cryptanalytic Attacks

We do not consider enumeration of all keys a valid cryptanalytic attack, since no well-designed cryptosystem is susceptible to such an approach. The types of legitimate attacks which we consider can be classified in three categories.

    1.   **Ciphertext-only Attack.**
    2.   **Known Plaintext Attack.**
    3.   **Chosen Plainext Attack.**

**Ciphertext-only Attack.**
The cryptanalyst intercepts one or more messages all encoded with the same encryption algorithm.

**Goal:** Recover the original plaintext or plaintexts, to discover the deciphering key or find an algorithm for deciphering subsequent messages enciphered with the same key.

**Known Plaintext Attack.**
The cryptanalyst has access to not only the ciphertext, but also the plaintext for one or more of the messages.

**Goal:** Recover the deciphering key or find an algorithm for deciphering subsequent messages (or the remaining plaintext) enciphered which use the same key.

**Chosen Plainext Attack.**
The cryptanalyst has access to ciphertext for which he or she specified he plaintext.

**Goal:** Recover the discover the deciphering key or find an algorithm for deciphering subsequent messages enciphered with the same key.

# Frequency Analysis

Given a sample of an English language newspaper text (stripped of spaces, punctuation and other extraneous characters) the following gives the approximate percentage of occurrences of each character.

| A | B | C | D | E | F | G | H | I | J | K | L | M |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 7.3 | 0.9 | 3.0 | 4.4 | 13 | 2.8 | 1.6 | 3.5 | 7.4 | 0.2 | 0.3 | 3.5 | 2.5 |

| N | O | P | Q | R | S | T | U | V | W | X | Y | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 7.8 | 7.4 | 2.7 | 0.3 | 7.7 | 6.3 | 9.3 | 2.7 | 1.3 | 1.6 | 0.5 | 1.9 | 0.1 |

The relative frequencies can change according to subject matter and style of writing, but it is still possible to pick out those characters with a high frequency of occurence and those which are rare:

**High frequency:** Vowels {E, I, O, A} and consonants {T, N, R, S}
**Low frequency:** {J, K, Q, X, Z}.

# Examples of Cryptanalysis

Let's begin with some ciphertext from a previous lecture.

```
QWMMPQDVKUVFDTXJQVDBOPIDUHDQQUGDLAMWJGXBGURRBPBURMKULDVX
OOKUJUOVDJQDGBWHLDJQQMUODQUBIMWBOVWUVXPBUBIOKUBGXBGURROK
UJUOVDJQVPWMMDJOUQDVKDBVKDCDAQXEDFKXOKLPWBIQVKDQDOWJXVAP
TVKDQAQVDHXQURMKULDVXOOKUJUOVDJQVKDJDTPJDVKDVPVURBWHLDJP
TCDAQXQPTDBPJHPWQQXEDBDNDJVKDRDQQFDFXRRQDDVKUVQXHMRDQWLQ
VXVWVXPBXQNDJAQWQODMVXLRDVPOJAMVUBURAVXOUVVUOCQ
```

We find the following character counts, scaled to that of a 1000 character input text.

| | | | | | |
|---|---|---|---|---|---|
| A | 24.0964 | J | 54.2169 | S | 0 |
| B | 54.2169 | K | 51.2048 | T | 15.0602 |
| C | 9.0361 | L | 24.0964 | U | 81.3253 |
| D | 129.5181 | M | 36.1446 | V | 99.3976 |
| E | 6.0241 | N | 6.0241 | W | 39.1566 |
| F | 12.0482 | O | 57.2289 | X | 57.2289 |
| G | 18.0723 | P | 45.1807 | Y | 0 |
| H | 18.0723 | Q | 96.3855 | Z | 0 |
| I | 12.0482 | R | 39.1566 | | |

The distributions look like a frequency preserving substitution cipher. We guess that the enciphering takes E $\mapsto$ D and T $\mapsto$ V or T $\mapsto$ Q. The most frequent characters are D, V, Q, V, U, O, J, K, B and E, N, S, Y, Z are of lowest frequency.

Equating high frequency and low frequency characters, we might first guess

$$\{\texttt{E}, \texttt{I}, \texttt{O}, \texttt{A}, \texttt{T}, \texttt{N}, \texttt{R}, \texttt{S}\} \mapsto \{\texttt{D}, \texttt{V}, \texttt{Q}, \texttt{U}, \texttt{O}, \texttt{J}, \texttt{K}, \texttt{B}\}$$

and

$$\{\texttt{J}, \texttt{K}, \texttt{Q}, \texttt{X}, \texttt{Z}\} \mapsto \{\texttt{E}, \texttt{N}, \texttt{S}, \texttt{Y}, \texttt{Z}\}$$

How would you go about reconstructing the entire text?

# Cryptanalysis: Frequency Analysis

**Index of Coincidence.**

In the 1920's William Friedman introduced the *index of coincidence* as a measure of the variation of character frequencies in text from a uniform distribution. The index of coincidence of a text space (e.g. that of all plaintext or ciphertext) is defined to be the probability that two randomly chosen characters are equal. In a given language over an alphabet of size $n$, suppose that $p_i$ is the probability of a random character is the $i$-th character. Then the index of coincidence in that language is:

$$\sum_{i=1}^{n} p_i^2.$$

Over an alphabet of 26 characters, the coincidence index of random text is

$$\sum_{i=1}^{26} (1/26)^2 = 1/26 \cong 0.0385.$$

For English text, the coincidence index is around 0.0661. For a finite string of length $N$, we determine the index of coincidence to be:

$$\frac{\sum_{i=1}^{n} n_i(n_i - 1)}{N(N-1)},$$

where $n_i$ is the number of occurrences of the $i$-th character in the string.

# Cryptanalysis: Recognizing Periodity

**Theorem 1** *The expected index of coincidence of a ciphertext of length $N$, output from a period $m$ cipher, which is defined by $m$ independent substitutions ciphers at each position in arithmetic progression the $i + jm$, is*

$$\frac{1}{m}\left(\frac{N-m}{N-1}\right) i_{\mathcal{L}} + \frac{(m-1)}{m}\left(\frac{N}{N-1}\right) i_n,$$

*where $i_{\mathcal{L}}$ is the index of coincidence of the language $\mathcal{L}$, the size of the alphabet is $n$, and and $i_n = 1/n$ is the index of coincidence of random text in that language.*

Consider the ciphertext enciphered with a Vigenère cipher:

```
OOEXQGHXINMFRTRIFSSMZRWLYOWTTWTJIWMOBEDAXHVHSFTRIQKMENXZ
PNQWMCVEJTWJTOHTJXWYIFPSVIWEMNUVWHMCXZTCLFSCVNDLWTENUHSY
KVCTGMGYXSYELVAVLTZRVHRUHAGICKIVAHORLWSUNLGZECLSSSWJLSKO
GWVDXHDECLBBMYWXHFAOVUVHLWCSYEVVWCJGGQFFVEOAZTQHLONXGAHO
GDTERUEQDIDLLWCMLGZJLOEJTVLZKZAWRIFISUEWWLIXKWNISKLQZHKH
WHLIEIKZORSOLSUCHAZAIQACIEPIKIELPWHWEUQSKELCDDSKZRYVNDLW
TMNKLWSIFMFVHAPAZLNSRVTEDEMYOTDLQUEIIMEWEBJWRXSYEVLTRVGJ
KHYISCYCPWTTOEWANHDPWHWEPIKKODLKIEYRPDKAIWSGINZKZASDSKTI
TZPDPSOILWIERRVUIQLLHFRZKZADKCKLLEEHJLAWWVDWHFALOEOQW
```

The coincidence index of this text is 0.0439. This would suggest a period of approximately 5. We will see that this is a bad estimate. Note that this doesn't disprove the theorem, it just shows that the statistical errors are too great and that we would need a much larger sample size to converge to this theoretical expectation, or that the substitutions employed were not independent.

**Exercise.** Explain why ciphertext for a particular key need not follow the behavior predicted by the theorem.

**Decimation of Sequences.**
For a sequence $S = s_1 s_2 s_3 \ldots$ and positive integers $m$ and $k$ such that $1 \le k \le m$, we denote the $k$-th decimation of period $m$ as the sequence $s_k s_{m+k} s_{2m+k} \ldots$. If $S$ is a ciphertext string (a sequence of characters in the alphabet $\mathcal{A}'$) enciphered by a cipher with period $m$, then the decimations of period $m$ capture the structure of the cipher without periods.

If we take the previous ciphertext and average the coincidence indices of each of the $k$-the decimated sequences of period $m$, we find:

| $m$ | CI | $m$ | CI | $m$ | CI |
|---|---|---|---|---|---|
| 1 | 0.0439 | 6 | 0.0424 | 11 | 0.0653 |
| 2 | 0.0438 | 7 | 0.0442 | 12 | 0.0408 |
| 3 | 0.0435 | 8 | 0.0414 | 13 | 0.0445 |
| 4 | 0.0434 | 9 | 0.0438 | 14 | 0.0418 |
| 5 | 0.0421 | 10 | 0.0407 | 15 | 0.0423 |

From this table, the correct period, 11, is obvious.

**Exercise.** What do you expect to see in such a table if the period is composite? Hint: consider the period of the decimated sequence, and apply the theorem.

**Kasiski method.**
The Prussian military officer Friedrich Kasiski made the following observation on the Vigenère cipher in 1863. If a frequently occurring pattern, such as THE is aligned at the same position with respect to the period, then the same three characters will appear in the ciphertext, at a distance which is an exact multiple of $m$.

By looking for frequently occurring strings in the ciphertext, and measuring the most

frequent divisors of the displacements of these strings, it is often possible to identify the period, hence to reduce to a monoalphabetic substitution.

To return to our sample ciphertext, we find that the three substrings SYE, ZKZ, and KZA each occur three times. The positions of these occurrences are:

$$\begin{array}{ll} \text{SYE:} & 122,\ 196,\ 383 \\ \text{ZKZ:} & 252,\ 439,\ 472 \\ \text{KZA:} & 253,\ 440,\ 473 \end{array}$$

Note that ZKZ and KZA are substrings of the four character string ZKZA appearing three times in the ciphertext! Moreover two of the occurences of the string SYE appear as substrings of the longer string SYEV.

Now we look for common divisors of the differences between the positions of the frequently occurring substrings.

$$\begin{array}{llll} 196 - 122 & = 2 \cdot 37 & 439 - 252 & = 11 \cdot 17 \\ 383 - 196 & = 11 \cdot 17 & 472 - 439 & = 3 \cdot 11 \\ 383 - 122 & = 3^2 \cdot 29 & 472 - 252 & = 2^2 \cdot 5 \cdot 11 \end{array}$$

We see that our guess of 11 for the period appears as a divisor of the distances between each of the occurrences of the common four character substring, and divides one of the differences of the other three string characters.

**Exercise.** If 11 is the correct period, why does it not appear in all of the differences above? Which of the occurences can be attributed as random?

## Breaking the Vigenère Cipher

Now that we have established that the period is 11, we can write the ciphertext in blocks and look at the strings which occur frequently at the same position within blocks.

|    | 1          | 2          | 3          | 4          |
|----|------------|------------|------------|------------|
| 1  | OOEXQGHXINM | FRTRIFSSMZR | WLYOWTTWTJI | WMOBEDAXHVH |
| 2  | SFTRIQKMENX | ZPNQWMCVEJT | WJTOHTJXWYI | FPSVIWEMNUV |
| 3  | WHMCXZTCLFS | CVNDLWTENUH | SYKVCTGMGYX | SYELVAVLTZR |
| 4  | VHRUHAGICKI | VAHORLWSUNL | GZECLSSSWJL | SKOGWVDXHDE |
| 5  | CLBBMYWXHFA | OVUVHLWCSYE | VVWCJGGQFFV | EOAZTQHLONX |
| 6  | GAHOGDTERUE | QDIDLLWCMLG | ZJLOEJTVLZK | ZAWRIFISUEW |
| 7  | WLIXKWNISKL | QZHKHWHLIEI | KZORSOLSUCH | AZAIQACIEPI |
| 8  | KIELPWHWEUQ | SKELCDDSKZR | YVNDLWTMNKL | WSIFMFVHAPA |
| 9  | ZLNSRVTEDEM | YOTDLQUEIIM | EWEBJWRXSYE | VLTRVGJKHYI |
| 10 | SCYCPWTTOEW | ANHDPWHWEPI | KKODLKIEYRP | DKAIWSGINZK |
| 11 | ZASDSKTITZP | DPSOILWIERR | VUIQLLHFRZK | ZADKCKLLEEH |
| 12 | JLAWWVDWHFA | LOEOQW      |            |            |

Some of the longer strings which appear more than one at distances which are a multiple of 11 are given in the following table. The first column indicates the number of times the full string appears.

| # | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|----|----|
| 3 | Z | A |   |   |   |   |   |   |   | Z  | K  |
| 2 |   |   | N | D | L | W | T |   |   |    |    |
| 2 |   |   |   | P | W | H | W | E |   |    |    |
| 2 | V |   |   |   |   |   |   |   | S | Y  | E  |
| 2 |   |   |   |   |   | L | W | C |   |    |    |

Now let's guess what the translations are at each of the periods. The following is a table of common characters in each of the 11 decimations of period 11, organized by the numbers of their appearances.

| $i$ | 9 | 8 | 7 | 6 | 5 |
|-----|---|---|---|---|---|
| 1 |   |   |   | S,W,Z | V |
| 2 |   |   |   | L | A |
| 3 |   |   |   | F | T |
| 4 |   |   | D | O | R |
| 5 |   |   | L |   | I,W |
| 6 | W |   |   |   | L |
| 7 | T |   |   | H | W |
| 8 |   |   |   | I,S,X | E |
| 9 |   |   | E |   | H |
| 10 |   |   | Z |   | E,Y |
| 11 |   |   | I |   |   |

These characters are not themselves the characters in the key, but if we assume that one of these frequently occuring characters is the image of E, then we can make a guess at the key. The table below gives the enciphering characters which take the corresponding character in the previous table to E.

| $i$ | 9 | 8 | 7 | 6 | 5 |
|-----|---|---|---|---|---|
| 1 |   |   |   | M,I,F | J |
| 2 |   |   |   | T | E |
| 3 |   |   |   | Z | L |
| 4 |   |   | B | Q | N |
| 5 |   |   | T |   | W,I |
| 6 | I |   |   |   | T |
| 7 | L |   |   | X | I |
| 8 |   |   |   | W,M,H | A |
| 9 |   |   | A |   | X |
| 10 |   |   | F |   | A,G |
| 11 |   |   | W |   |   |

Checking possible keys, the partial key I****IL*A*W gives the following text which is suggestive of English:

```
            1           2           3           4
  1   W****OS*I*I   N****ND*M*N   E****BE*T*E   E****LL*H*D
  2   A****YV*E*T   H****UN*E*P   E****BU*W*E   N****EP*N*R
  3   E****HE*L*O   K****EE*N*D   A****BR*G*T   A****IG*T*N
  4   D****IR*C*E   D****TH*U*H   O****AD*W*H   A****DO*H*A
  5   K****GH*H*W   W****TH*S*A   D****OR*F*R   M****YS*O*T
  6   O****LE*R*A   Y****TH*M*C   H****RE*L*G   H****NT*U*S
  7   E****EY*S*H   Y****ES*I*E   S****WW*U*D   I****IN*E*E
  8   S****ES*E*M   A****LO*K*N   G****EE*N*H   E****NG*A*W
  9   H****DE*D*I   G****YF*I*I   M****EC*S*A   D****OU*H*E
 10   A****EE*O*S   I****ES*E*E   S****ST*Y*L   L****AR*N*G
 11   H****SE*T*L   L****TH*E*N   D****TS*R*G   H****SW*E*D
 12   R****DO*H*W   T****E
```

**Exercise.** Complete the deciphering. What do you note about the relation between the text and the enciphering or deciphering key?