

Code Breaking

So far we have focused on Vigenère ciphers, and their reduction to monoalphabetic substitutions. Here we show how to use **Magma** to complete the final step of breaking these ciphers. Recall that the reduction to monoalphabetic substitution is done by the process of *decimation*, by which we lose all 2-character frequency structure of the language. A more sophisticated approach will be necessary for breaking more complex ciphers.

Correlation. We first introduce the concept of correlation of two functions. Let $X = x_1, x_2, \dots, x_n$ and $Y = y_1, y_2, \dots, y_n$ be two finite sequences of real numbers, each of length n . We define the correlation of the two sequences to be

$$\text{Corr}(X, Y) = \frac{\sum_{i=1}^n (x_i - \mu_X)(y_i - \mu_Y)}{\sigma(X)\sigma(Y)}$$

and where μ_X and μ_Y are the respective *means* of the sequences X and Y :

$$\mu_X = \frac{1}{n} \sum_{i=1}^n x_i, \quad \mu_Y = \frac{1}{n} \sum_{i=1}^n y_i,$$

and the terms in the denominators are:

$$\sigma(X) = \left(\sum_{i=1}^n (x_i - \mu_X)^2 \right)^{1/2}, \quad \sigma(Y) = \left(\sum_{i=1}^n (y_i - \mu_Y)^2 \right)^{1/2},$$

called the *standard deviations* of X and Y .

The correlation of two sequences will be a real number between 1 and -1 , which measures the linear relation between two sequences. On the following page we give a **Magma** function which computes the sequence of correlations of each of the cyclic translations of two sequences.

- 1. Correlations of sequence translations.** The following code from the course cryptography package finds the correlations between the affine translations of two sequences. Do you understand the syntax? Ask or refer to the online **Magma** handbook where necessary.

```
function TranslationCorrelations(S1,S2)
// The sequence of correlations of the sequence S1 with the
// cyclic translations of the sequence S2.
```

```

n := #S1;
error if n ne #S2, "Arguments must be of the same length.";

// Compute the mean value of each sequence:
mu1 := &+[ S1[k] : k in [1..n] ]/n;
mu2 := &+[ S2[k] : k in [1..n] ]/n;

// Compute the standard deviations of each sequence:
sig1 := Sqrt(&+[ (S1[k]-mu1)^2 : k in [1..n] ]);
sig2 := Sqrt(&+[ (S2[k]-mu2)^2 : k in [1..n] ]);

sig := sig1*sig2;
CorrSeq := [ ];
for j in [1..n] do
    Corr := &+[ (S1[i] - mu1) * (S2[ij] - mu2) / sig
                where ij := ((i+j-1) mod n) + 1 : i in [1..n] ];
    Append(~CorrSeq,<j,Corr>);
end for;
return CorrSeq;
end function;

```

Solution Given two frequency distributions, S_1 and S_2 this function computes the correlation of S_1 with each of the cyclic shifts of S_2 . The function will be called by **TranslationMatches** in the exercise below.

- 2. Breaking Vigenère ciphers.** A Vigenère cipher is reduced to an translation cipher by the process of decimation. How does the above function solve the problem of finding the affine translation?

For completeness we give a function, using the previous one, which matches the frequency distribution of a string with a given standard distribution. This uses the function `FrequencyDistribution`, which is part of the course cryptography package.

```

function TranslationMatches(S,F,r)
    // INPUT:
    // S : Test string.
    // F : Sequence of standard frequencies for the language.
    // r : A real number between 0 and 1.
    // OUTPUT:
    // Returns integers k such that affine translation
    // of S by k has correlation at least r with the standard
    // frequencies given by the real sequence F.

    X := FrequencyDistribution(S);
    CorrSeq := TranslationCorrelations(X,F);

```

```

    return [ x[1] : x in CorrSeq | x[2] gt r ];
end function;

```

Use this function on the Vigenère ciphertext samples from the web page the break the enciphering. Recall that you will have to use the functions `Decimation` and `CoincidenceIndex` to first reduce a Vigenère cipher to a monoalphabetic one.

Solution We have already surmised that the first sample ciphertext, `cipher01.txt`, is output from a Vigenère cipher of period 11. Using the above function, we determine candidate translations:

```

> PT := StripEncoding(Read("blackcat.txt"));
> FD := FrequencyDistribution(PT);
> CT := StripEncoding(Read("cipher01.txt"));
> for i in [1..11] do
>   Ci := Decimation(CT,i,11);
>   TranslationMatches(Ci,FD,0.50);
> end for;
[ 8 ]
[ 19 ]
[ 26 ]
[ 16 ]
[ 22 ]
[ 8 ]
[ 11 ]
[ 22 ]
[ 26 ]
[ 9 ]
[ 22 ]

```

A translation by 8 corresponds to the 9th character, I, by 19 to the 20th character, T, etc., giving the key `ITAQWILWAJW`. Enciphering with respect to this key gives the plaintext: `WHENMOSTIWINK...` Note that the inverse key is `SHAKESPEARE` – why is this a bad choice of key?

3. Breaking substitution ciphers. Suppose that rather than an affine translation, you have reduced to an arbitrary simple substitution. We need to undo an arbitrary permutation of the alphabet. For this purpose we define maps into Euclidean space:

a. $\mathcal{A} \rightarrow \mathcal{A}^2 \rightarrow \mathbb{R}^2$ defined by

$$X \mapsto XX \mapsto (P(X), P(XX)).$$

b. $\mathcal{A} \rightarrow \mathcal{A}^2 \rightarrow \mathbb{R}^3$ defined by

$$X \mapsto XY \mapsto (P(X), P(XY|Y), P(YX|Y))$$

for some fixed character Y .

See the document `digraph_frequencies.pdf` for standard vectors for the English language.

Solution These maps can be applied to the solution of substitution ciphers by finding nearest elements to a known standard for the English language. For instance, assume that the ciphertext image X of \mathbf{E} has been identified, one can look for pairs XY which are the image of the plaintext pair \mathbf{ER} , by searching for a nearest vector to:

$$(0.05674, 0.13071, 0.11352).$$

This will determine the ciphertext image Y of \mathbf{R} . This bootstrapping procedure successively determines the substitution from the digraph frequencies of the ciphertext.

- 4. Breaking transposition ciphers.** In order to break transposition ciphers it is necessary to find the period m , of the cipher, and then to identify positions i and j within each block $1 + km \leq i, j \leq (k + 1)m$ which were adjacent prior to the permutation of positions. Suppose we guess that m is the correct period. Then for a ciphertext sample $C = c_1c_2\dots$, and a choice of $1 \leq i < j \leq m$, we can form the digraph decimation sequence $c_i c_j, c_{i+m} c_{j+m}, c_{i+2m} c_{j+2m}, \dots$

Two statistical measures that we can use on ciphertext to determine if a digraph sequence is typical of the English language are a digraph *coincidence index*

$$\sum_{X \in \mathcal{A}} \sum_{Y \in \mathcal{A}} \frac{n_{XY}(n_{XY} - 1)}{N(N - 1)}$$

where N is the total number of character pairs, and n_{XY} is the number of occurrences of the pair XY , and the *coincidence discriminant*:

$$\sum_{X \in \mathcal{A}} \sum_{Y \in \mathcal{A}} \left(\frac{n_{XY}}{N} - \sum_{Z \in \mathcal{A}} \frac{n_{XZ}}{N} \sum_{Z \in \mathcal{A}} \frac{n_{ZY}}{N} \right)^2.$$

The first term is the frequency of XY , and the latter is the product over the frequencies of X as a first character and Y as a second character. The coincidence discriminant measures the discrepancy between the probability space of pairs XY and the product probability space.

What behavior do you expect for the coincidence index and coincidence discriminant of the above digraph decimation, if i and j were the positions of originally adjacent characters? Test your hypotheses with decimations of “real” English text.

```
function CoincidenceDiscriminant(S)
    // INPUT: A sequence of 2-character strings, produced
    // as decimation of transposition ciphertext, or of
    // adjacent characters in some sample plaintext.
    // OUTPUT: A measure of the difference of probability
    // of association of two characters, relative to their
    // independent probabilities.
```

```

C2S := CodeToString;
AA := [ C2S(64+i)*C2S(64+j) : i, j in [1..26] ];
FD1 := FrequencyDistribution(&*[ s[1] : s in S ]);
FD2 := FrequencyDistribution(&*[ s[2] : s in S ]);
N := #S;
F2D := [ RealField() | 0 : i in [1..26^2] ];
for s in S do
    F2D[Index(AA,s)] += 1/N;
end for;
return &+[ (F2D[i+26*(j-1)]-FD1[i]*FD2[j])^2 : i,j in [1..26] ];
end function;

```

Why can we assume that $i < j$ in the digraph sequence? What is the obstacle to extending these statistical measures from two to more characters?

Solution If i and j are the ciphertext images of adjacent positions $k, k + 1$, in each block of length m , then the sequence

$$c_i c_j, c_{i+m} c_{j+m}, c_{i+2m} c_{j+2m}, \dots$$

will have the coincidence index and coincidence discriminant of the plaintext. Note that the measures are invariant under a substitution, so can be used to break a combination substitution-transposition cipher, by first breaking the transposition. The result will be a sequence i_1, i_2, \dots, i_m of indices of positions which “want” to be associated.

Note that the measures, coincidence index and coincidence discriminant, will also be the same for the sequence

$$c_j c_i, c_{j+m} c_{i+m}, c_{j+2m} c_{i+2m}, \dots$$

so we do not directly distinguish the correct order from its reverse, i_m, \dots, i_2, i_1 . This one bit of information can be determined at the end, with the savings of being able to assume $i < j$ in testing for statistically associated pairs $\{i, j\}$.

Note also that for the incorrect period m there will be little or no tendency for statistical association of characters, so by first varying the triples (i, j, m) , with fixed $i = 1$ and $1 < j \leq m$, we can determine the probable period m and then recover the entire sequence i_1, i_2, \dots, i_m by letting i vary.