# The Vigenère cipher.

This is a more elaborate simple substitution cipher, than the shift cipher (which we discussed before).

In a way, the Vigenère cipher is a generalization of the shift ciphers, with the shifts being more 'chaotic'.

The cipher is performed as follows:

(1) choose a keyword (responsible for the shifts);

(2) shift $k^{th}$ letter of the plaintext by the number corresponding to $k \pmod{\ell}$'s letter of the keyword ($\ell$ = length (keyword)).

Example. Consider the plaintext 'bitcoin' and Keyword 'cat' = '3 1 20' (using the correspondence $A \leftrightarrow 1$, $B \leftrightarrow 2$, ---, $Z \leftrightarrow 26$).

$$b \xrightarrow{+3} e$$
$$i \xrightarrow{+1} y$$
$$t \xrightarrow{+20} n$$
$$c \xrightarrow{+3} f$$
$$o \xrightarrow{+1} p$$
$$i \xrightarrow{+20} c$$
$$n \xrightarrow{+3} k$$

The ciphertext is 'ejnfpck'.

# Cryptanalysis of the Vigenère cipher.

At some point Vigenère-type ciphers were thought to be 'unbreakable', but this is far from being true. Our next goal is to get acquainted with the basic statistical tools in cryptanalysis.

**Def-n.** Let $s = a_1 a_2 \ldots a_n$ be a text (string of letters). The index of coincidence of $s$, denoted via $IndCo(s)$, is the probability that two randomly chosen characters coincide.

**Example.** Let's take $s = aabbcc$. There are $\binom{6}{2} = \frac{6!}{2!4!} = 15$ pairs and 3 of them consist of coinciding letters ('aa', 'bb' and 'cc'). Hence, $IndCo(s) = \frac{3}{15} = 0.2$.

Next we obtain a general formula. For this purpose, denote the number of occurrences of the $i^{th}$ letter of the alphabet in $s$ by $F_i$, then

$$IndCo(s) = \frac{\sum_{i=1}^{26} \binom{F_i}{2}}{\binom{n}{2}} = \frac{\sum_{i=1}^{26} F_i(F_i - 1)}{n(n-1)},$$

<u>Remark.</u> If the string $s$ consists of random characters (the occurence of each letter in $s$ is equally likely), then

$$IndCo(s) = \frac{26 \cdot \frac{n}{26}\left(\frac{n}{26}-1\right)}{n(n-1)} = \frac{\left(\frac{n}{26}-1\right)}{n-1} \xrightarrow[n \to \infty]{} \frac{1}{26} \approx 0.0385.$$

$$\left(F_i = \frac{n}{26}\right)$$

However, if $s$ consists of an actual text (in English), then $IndCo(s) \approx 0.0685$ (length($n$) sufficiently large). Notice that this value is almost 2 times greater than for a random text!

<u>Another important remark.</u> $IndCo(s)$ does not change (is invariant) under permutation letters, i.e. withstands any simple substitution.

<u>Strategy to break the Vigenère cipher.</u>

<u>Step 1.</u> Find the length of the keyword ($\ell$).

Let's take a number $k$ and test if $k = \ell$. In order to perform such a test, consider for every $1 \le i \le k$ the subtext $S_i$ of $s$, where

$$S_i = a_i \, a_{i+k} \, a_{i+2k} \dots$$

$$(s = a_1 a_2 \dots a_n) \text{ is the ciphertext}$$

If our guess was correct $(K=\ell)$, then $IndCo(s_i)$ will be close to the one for a text in English (see the remark above). On the other hand, if $K \neq \ell$, then $IndCo(s_i)$ will be close to the one for the random text.

To sum up, we take the average of the indices of coincidence and compare it to $0.0385$ and $0.0685$:

- if $\dfrac{\sum_{i=1}^{k} IndCo(s_i)}{k} \sim 0.0685$, then $K=\ell$ is likely.

- if $\dfrac{\sum_{i=1}^{k} IndCo(s_i)}{k} \sim 0.0385$, then $K=\ell$ is unlikely.

So we try $K=1, 2, 3$, etc. and, as the keyword has some 'adequate' length, at some point find $\ell$.

Step 2. Next one needs to figure out the actual keyword.

Def-n. Let $s = a_1 \ldots a_n$ and $t = b_1 \ldots b_m$ be two strings of letters. The mutual index of coincidence of $s$ and $t$ is the probability that a randomly chosen symbol in $s$ and a randomly chosen symbol in $t$

will be the same.

$$\text{MutIndCo}(s,t) := \frac{1}{nm} \sum_{i=1}^{26} F_i(s) F_i(t),$$

where $F_i(s)$ is the number of occurrences of the $i^{th}$ letter in $s$ and $F_i(t)$ in $t$.

Example. Let $s$ = 'Catalonia' and $t$ = 'Barcelona'.

Then $n$ = len$(s)$ = 9, $m$ = len$(t)$ = 9 and we find

$$\text{MutIndCo}(s,t) = \frac{1}{9\cdot 9}\underset{\substack{\qquad\;\; 'a'\;\;\;\; 'c'\;\; 'l'\;\; 'n'\;\; 'o'}}{(3\cdot 2 + 1\cdot 1 + 1\cdot 1 + 1\cdot 1 + 1\cdot 1)} = \frac{10}{81}.$$

So, they have something incommon.

(✬)Remark. If two strings $s$ and $t$ are encryptions of plaintext with the help of the same SSC, then the value of $\text{MutIndCo}(s,t)$ will be larger. This index is an analogue of correlation in probability (that's my intuition at least).

Next we will use the MutIndCo to 'finish decrypting' Vigenere cipher.

Let
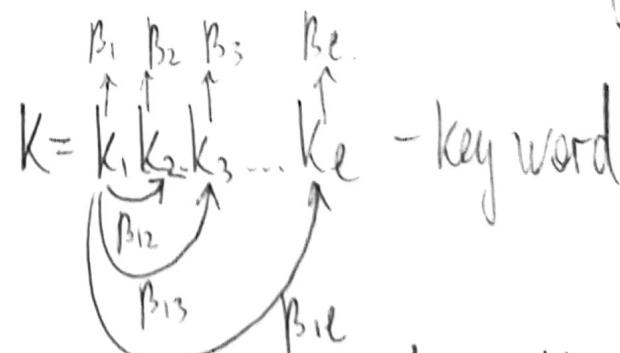$$S^{(1)} = a_1 a_{1+\ell} a_{1+2\ell} \cdots$$
$$S^{(2)} = a_2 a_{2+\ell} a_{2+2\ell} \cdots$$
$$S^{(\ell)} = a_\ell a_{2\ell} a_{3\ell} \cdots$$

$P = p_1 p_2 p_3 \cdots$ plaintext
$S = a_1 a_2 a_3 \cdots$ ciphertext

$$\overset{\beta_1\;\;\beta_2\;\;\beta_3\qquad\;\; \beta_\ell}{\underset{}{K = k_1 k_2 k_3 \ldots k_\ell}} - \text{key word}$$

$\beta_{12}$
$\beta_{13}$
$\beta_{1\ell}$

$\beta_i$ = index of letter $k_i$ in the alphabet

$\beta_{ii'} := \beta_i - \beta_{i'}$

<u>Remark.</u> Notice that $s^{(i)} = a_i a_{i+\ell} a_{i+2\ell} \ldots$ is encryption of $p^{(i)} = p_i p_{i+\ell} p_{i+2\ell} \ldots$ via shift cipher (the shift is by $k_i$ with index $\beta_i$).

Next we compare the mutual indices of coincidence $\mathrm{MutIndCo}(s^{(1)}, s^{(i)}$ shifted by $h_i)$ for different values of $h_i$ between $0$ and $25$. As for $h_i = \beta_{1i}$ we get that both $s^{(1)}$ and $s^{(i)}$ are shifts of the corresponding parts of plaintext by $\beta_1$ (notice that $\beta_i + \beta_1 - \beta_i = \beta_1$), the correspon-ding index $\mathrm{MutIndCo}(s^{(1)}, s^{(i)}$ shifted by $\beta_{1i})$ will be the largest (see Remark $(\#)$ on the previous page). This allows to find the values of $\beta_{12}, \beta_{13}, \ldots, \beta_{1\ell}$. Now it suffices to establish $\beta_1$ in order to decide the key word $K$. For this we simply try the $26$ possible options.

<u>Remark.</u> The $\beta_{1i}$'s are found independently, so the complexity is $\vartheta(\ell)$, not $\vartheta(25^{\ell-1})$ as for a nested loop.