

Cryptanalysis of Vigenère cipher

A alphabet, Latin

P_1, P_2, \dots, P_t on A , t block size (period)

$$m = \begin{array}{c|c|c|c|c} m_1 & \dots & m_t & | & m_{t+1} \dots \\ \downarrow & & \downarrow & & \downarrow \\ P_1(m_1) & \dots & P_t(m_t) & | & P_1(m_{t+1}) \dots \end{array}$$

P_1, \dots, P_t are shifts on A

random

simple Vigenère

full Vigenère

Problem: given cipher-text recover pl.-text.

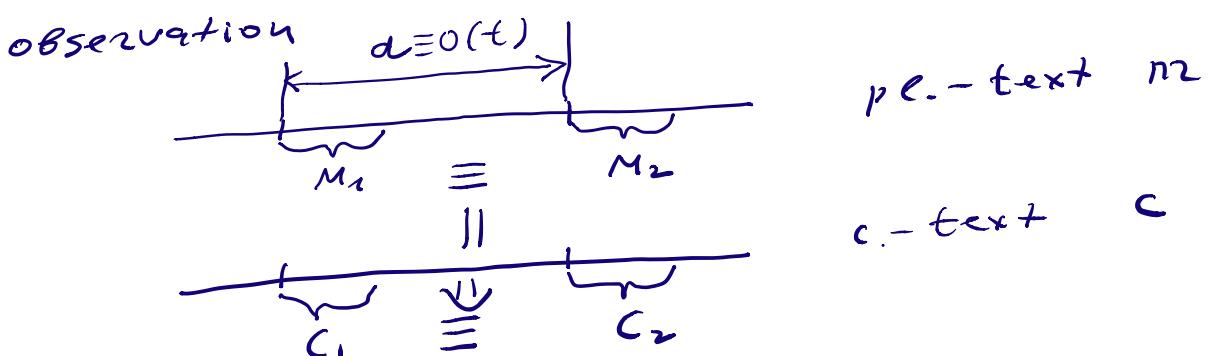
we do not know t .

1) find t

2) use frequency analysis to find m .

How to find t ?

1. Kasiski test. $\Leftrightarrow t/d$ is a multiple of t



Search c.-text for pairs of identical segments and note distance between starting points of them

collect distances d_1, \dots, d_s

very likely $t = \gcd(d_1, \dots, d_s)$

some distance $d_{s+1} \not\equiv 0 \pmod{t}$

$\gcd(d_1, \dots, d_s, d_{s+1})$ very small like 1, 2, ...

2. Autocorrelation (more Randy)

	$c = c_1 c_2 \dots c_N$	$\# \text{ occurrences } c_i = c_{i+d}$	$c_1 c_2 \dots c_N$
distanced		$\frac{1}{N-d}$	$c_1 c_2 \dots c_N$
1	≈ 0.038	$\frac{1}{26}$	$c_1 c_2 \dots c_N$
2	≈ 0.038	:	$c_1 c_2 \dots c_N$
$t-1$	≈ 0.038		
t		first relatively large value $\approx 0.065 = \sum \alpha^2$	

recover t .

Why?

Assume pl.-text is randomly chosen.

$$i \text{ fixed} \quad P_r(c_i = c_{i+d}) = P_r(\underbrace{P_u(m_i)}_{\sim c_i} = \underbrace{P_v(m_{i+d})}_{\sim c_{i+d}}) =$$

$$\boxed{u \equiv i \pmod t \\ v \equiv i+d \pmod t}$$

$$= P_r(m_i = \underbrace{\overline{P_u^{-1}}}_{R} \underbrace{P_v}_{\sim}(m_{i+d})) =$$

$$= \sum_{d \in A} P_r(m_i = R(m_{i+d}), m_{i+d} = d)$$

complete probability formula

$$= \sum_{d} P_r(m_i = R(d), m_{i+d} = d) \approx$$

joint even

$$\approx \sum_{d} P_r(m_i = R(d)) \cdot P_r(m_{i+d} = d)$$

by independence of m_i, m_{i+d}
for rel. large d .

$$= \sum_{d} q_{R(d)} \cdot q_d$$

1) $d \equiv 0 \pmod t \Rightarrow \boxed{u \equiv v \pmod t \Rightarrow P_u = P_v \Rightarrow R=1}$

$$\Pr(c_i = c_{i+a}) = \sum_{\alpha} q_{\alpha} \approx 0.065.$$

2) $d \not\equiv 0 \pmod{t} \Rightarrow R$ looks as a random subst.
average of $\sum_{\alpha} q_{R(\alpha)} \cdot q_{\alpha}$ over all subs. R .

$$\frac{1}{26!} \sum_R \sum_{\alpha} q_{R(\alpha)} \cdot q_{\alpha} = \\ = \frac{1}{26!} \sum_{\alpha} q_{\alpha} \sum_R q_{R(\alpha)} = \frac{25!}{26!} \sum_{\alpha} q_{\alpha} = \frac{1}{26} \approx 0.038$$

$$\sum_R q_{R(\alpha)} = \sum_{\beta \in A} \sum_{R: R(\alpha) = \beta} q_{\beta} = \\ = \sum_{\beta} q_{\beta} \sum_{R: R(\alpha) = \beta} 1 = 25! \sum_{\beta} q_{\beta} = 25! \\ \text{--- } \alpha - 25 \quad 25! \\ \text{--- } \beta - -$$

We now know t . How to find
pl-text?

$c_1 c_2 \dots c_N$

Write cipher-text as an array $t \times \approx N/t$
 $\approx N/t + P_t$

$C(1)$	$c_1 c_{1+t} c_{1+2t} \dots$	\leftarrow	$m_1 m_{1+t} m_{1+2t} \dots$
$C(2)$	$c_2 c_{2+t} c_{2+2t} \dots$	$\leftarrow P_2$	$m_2 m_{2+t} m_{2+2t} \dots$
\vdots			
$C(t)$	$c_t c_{2t} c_{3t} \dots$	$\leftarrow P_t$	$m_t m_{2t} m_{3t} \dots$

look at $C(i)$

Find frequencies of ch. in $C(i)$

$$\gamma_0 \gamma_1 \dots \gamma_{25} = \gamma_a \gamma_b \dots \gamma_z$$

$\gamma_{P_i(d)} \approx q_d$ frequency of $d \in A$ in long English texts.

- 1 full Vigenère $P_1 \dots P_t$ are random subs. out.
 guess on encryption of E, T, A, \dots
 do that for all $C(i)$, partly decrypt
 cipher-text and check the guesses.
 But if N is large enough, then it works.
 $N \approx 100 \cdot t$

2 simple Vigenère

$$P_i(d) = d + k_i \bmod 26$$

$$\begin{matrix} A & B & C & \dots & Z \\ 0 & 1 & 2 & \dots & 25 \end{matrix}$$

take row $C(i)$, compute frequencies in $C(i)$

$$\gamma_0 \gamma_1 \dots \gamma_{25}$$

$$\gamma_{d+k_i \bmod 26} \approx q_d$$

measure for possible shifts s :

$$M_s = \sum_d q_d \cdot \gamma_{d+s \bmod 26}$$

two cases

1) $s = k_i$ correct shift

$$M_s = \sum_d q_d \cdot \gamma_{d+k_i} \approx \sum_d q_d \cdot q_d \approx 0.065$$

2) $s \neq k_i$
 average of M_s (over all shifts s)
 $\approx \frac{1}{26} = 0.038$

correct shift is recovered by computing

26 values M_5 .

complexity of the attack.

$$t \left(\frac{N}{t} + 26^2 \right)$$

\uparrow
find frequencies
in $C(i)$

\uparrow
compute 26 of M_5
sum 26 terms

$$= N + t \cdot 26^2$$

some operations

compared with P_1, P_2, \dots, P_t brute force attack

26^t trials.

conclude efficient cryptanalysis.

Example. Lecture Notes

cipher-text

KPBpw 9zxrg ---
of length 332 characters.

frequencies	l	0.081
	p	0.078
	r	0.069
	x	0.066

0.137 -

assume a Vigenère cipher

find block size t

apply Kasiski test

7-cB. segments "aahjeyx" in positions

309, 209 in the c-text.

distance $309 - 209 = 100$

3-cR. "psx" in positions 23, 73
distance = $258 - 73 = 185$

$$\gcd(100, 185) = 5$$

there are a lot of 3-cR. identical segments in the cipher-text at distances multiple of 5

$$\Rightarrow t = 5.$$

array of 5 rows and $\frac{332}{5} \approx 66$

$C(1) = kgzzf\ldots$ of length 67

frequencies in $C(1)$

$$z = (0, 0, 0, \overset{\alpha}{\underset{0}{3}}, \overset{\beta}{\underset{1}{2}}, \overset{\gamma}{\underset{2}{1}}, \overset{\delta}{\underset{3}{6}}, \overset{\epsilon}{\underset{4}{7}}, \overset{\zeta}{\underset{5}{12}}, \underset{6}{\underline{\underline{7}}}, \ldots)$$

compute $(M_0, M_1, \ldots, M_{25}) =$

$$(0.032, 0.031, 0.041, \underset{0}{\frac{0.065}{2}}, \underset{1}{\underline{\underline{0.043}}}, \underset{2}{\approx \sum_x q_x^2}, \ldots)$$

$$\Rightarrow k_1 = 3$$

similar for $C(2), \ldots, C(5)$

recover cipher key $(3, 11, 4, 7, 19)$

decrypt cipher-text to

HE DID NOT KNOW....

Index of coincidence.

$x = x_1 x_2 \ldots x_N$ N-ch. string over

Latin alphabet A.

index $I_c(x)$ measures how close x to a sensible English text.

$T(x)$ probability that

definition: $I_c(x)$,
 two random x_i, x_j from x ($i \neq j$)
 are identical $x_i = x_j$.

introduce indicator

$$\gamma_{ij} = \begin{cases} 1 & x_i = x_j \\ 0 & x_i \neq x_j \end{cases} \quad 1 \leq i < j \leq N$$

$$I_c(x) = \frac{1}{\binom{N}{2}} \cdot \sum_{1 \leq i < j \leq N} \gamma_{ij}$$

$$\binom{N}{2} = \frac{N(N-1)}{2}$$

find expectation of $I_c(x)$

- 1) x random string of characters
- 2) $x = m$ random English text.

then we prove a Baudy formula to
 compute $I_c(x)$.

1) x random N -character string.

$$E I_c(x) = \frac{1}{\binom{N}{2}} \sum_{1 \leq i < j \leq N} E \gamma_{ij}$$

\mathbb{E} expectation is a linear operator

$$E \gamma_{ij} = 1 \cdot P_2(x_i = x_j) + 0 \cdot P_2(x_i \neq x_j) = \\ = P_2(x_i = x_j)$$

$$\Rightarrow E I_c(x) = \frac{1}{\binom{N}{2}} \sum_{1 \leq i < j \leq N} P_2(x_i = x_j) = *$$

$$\Pr(x_i = x_j) = \frac{1}{26}$$

$$* = \frac{1}{\binom{N}{2}} \cdot \binom{N}{2} \cdot \frac{1}{26} = \frac{1}{26} \approx \underline{\underline{0.038}}$$

2) $x = m$ random English text

$$EI_c(m) = \frac{1}{\binom{N}{2}} \sum_{1 \leq i < j \leq N} \Pr(x_i = m_i)$$

$$\Pr(x_i = m_i) = \sum_{\alpha \in A} \Pr(x_i = m_i, x_j = \alpha) =$$

$$= \sum_{\alpha} \Pr(x_i = \alpha, x_j = \alpha) = \underset{j=i+1 \dots m_i, m_{i+1}}{\text{TH}}$$

$$\approx \sum_{\alpha} \Pr(x_i = \alpha) \cdot \Pr(x_j = \alpha) = \sum_{\alpha} q_{\alpha}^2 \approx 0.065$$

$$\Rightarrow EI_c(m) \approx \underline{\underline{0.065}}$$

Why do we need $I_c(x)$?

automate decryption.

formula to compute $I_c(x)$

$$x = x_1 \dots x_N$$

$f_{\alpha} = \# \alpha \text{ appears in } x$.

$$I_c(x) = \frac{1}{\binom{N}{2}} \left(\sum_{\alpha \in A} \binom{f_{\alpha}}{2} \right) = \frac{\sum_{\alpha} f_{\alpha}(f_{\alpha}-1)}{N(N-1)}$$

$\# x_i = x_j$
 $i \neq j$

apply to text
 $x = \text{THE DATA ENCRYPTION STANDARD}$

α	P_α
A, T	4
D, N	3
E, R	2
C, H, I, O, S, Y	1

$$\Rightarrow I_c(x) =$$

$$= \frac{4 \cdot 3 + 4 \cdot 3 + 3 \cdot 2 + 3 \cdot 2 + 2 \cdot 1 + 2 \cdot 1 + 1 \cdot 0}{25 \cdot 24} =$$

$$= \frac{40}{600} \approx 0.066 \approx \sum q_\alpha^2$$

$\Rightarrow x$ is likely to be sensible English text.

index may be defined for digrams, trigrams, ...