

3.3. Estimating the key length of a Vigenère cipher

Stochastic model :

$$\mathcal{X} = \{0, \dots, m-1\} \quad \text{Alphabet}$$

k keyword length, n message length, $k \ll n$

$$M = (M_1, \dots, M_k, M_{k+1}, \dots, M_{2k}, M_{2k+1}, \dots, M_n)$$

$$\oplus K = (K_1, \dots, K_k, K_1, \dots, K_k, K_1, \dots, K_k)$$

$$C = (C_1, \dots, C_k, C_{k+1}, \dots, C_{2k}, C_{2k+1}, \dots, C_n)$$

$$M_i \text{ i.i.d.}, P(M_i = e) = p_e \quad (\text{known})$$

$$K_i \text{ i.i.d.}, P(K_i = e) = \frac{1}{m}$$

I_c : index of coincidence,

$$I_c = \frac{1}{\binom{n}{2}} \sum_{1 \leq i < j \leq n} Y_{ij}, \quad Y_{ij} = \begin{cases} 1, & C_i = C_j \\ 0, & \text{otherwise} \end{cases}$$

$$K_M = \sum_{e=0}^{m-1} p_e^2$$

Lemma 3.5. $E(I_c) = \frac{1}{k(n-1)} \left[(n-k) K_M + k(k-1) \frac{1}{m} \right] (*)$

Outline of the proof.

Consider 2 ~~or~~ cases:

$$1.) \quad i \equiv j \pmod{k}$$

$$E(Y_{ij}) = \sum_{l=0}^{m-1} p e^l = k_M$$

$$2.) \quad i \not\equiv j \pmod{k}$$

$$E(Y_{ij}) = \frac{1}{m}$$

$$\text{Finally: } E(\bar{I}_c) = \frac{1}{\binom{m}{2}} \sum_{i < j} E(Y_{ij})$$

$$= \frac{1}{\binom{m}{2}} \left[\sum_{\substack{i < j \\ i \equiv j}} E(Y_{ij}) + \sum_{\substack{i < j \\ i \not\equiv j}} E(Y_{ij}) \right]$$

$$= (*)$$

We are interested in k . Solve (*) for k :

$$(m-1) E(\bar{I}_c) = \frac{1}{k} \left(m \left(k_M - \frac{1}{m} \right) - \left(k_M - \frac{m}{m} \right) \right)$$

$$k = \frac{m \left(k_M - \frac{1}{m} \right)}{(m-1) E(\bar{I}_c) + k_M - \frac{m}{m}}$$

Application: Estimate $E(\bar{I}_c)$ by \bar{I}_c

$$\bar{I}_c = \frac{1}{n(n-1)} \sum_{e=1}^{n-1} n e (n e - 1)$$

By Lemma 3.3. Pf. : $\bar{I}_c \rightarrow E(\bar{I}_c) \quad (n \rightarrow \infty) \text{ a.e.}$

In German: $K_M = 0.0762, n = 26$

Hence:

$$\hat{k} = \frac{0.0377n}{(n-1)\bar{I}_c - 0.0385n + 0.0762}$$

If k is known, write C as follows

$$\hat{C} = \begin{pmatrix} c_1 & \dots & c_k \\ c_{k+1} & \dots & c_{2k} \\ \vdots & & \vdots \\ c_{sk+1} & \dots & c_n \end{pmatrix}$$

The columns are monoalphabetic, apply frequency analysis to the columns.

3.4. Vigenère cipher with running key

$$\begin{array}{cccc}
 & a_1 & a_2 & \dots & a_n \\
 \oplus & s_1 & s_2 & \dots & s_n \\
 \hline
 & c_1 & c_2 & \dots & c_n
 \end{array}
 \quad \text{(taken from a book)}$$

Frequency attack is possible, if (s_1, \dots, s_n) is from a nat. language.

Model: M_i r.v.s occurrence of plaintext char } stoch. ind.
 K_i r.v.s " " " " " key char.

Consider most frequ. char.: E, T, A, O, I, N, S (57%)

$$P((M_i, K_i) \in \{E, \dots, S\}^2) = 0.57^2 = 0.3249$$

About $\frac{1}{3}$ of all ciphertext char. are obtained by 'adding' 2 of the most frequent char.

Most frequent characters: (total 57.29%)

E - 12.51%, T - 9.25%, A - 8.04%, O - 7.60%, I - 7.26%, N - 7.09%, S - 6.54%

Example:

I T I S N I C E T O L E A R N A B O U T C R Y P T O G R A P H Y
E T A I S A N E L E M E N T O F T H E G R E E K A L P H A B E T
M M I A F I P I E S X I N K B F U V Y Z T V C Z T Z V Y A Q L R

Investigating the first five ciphertext characters:

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
M L K J I H G F E D C B A Z Y X W V U T S R Q P O N
M M

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
I H G F E D C B A Z Y X W V U T S R Q P O N M L K J
I I

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
A Z Y X W V U T S R Q P O N M L K J I H G F E D C B
A A

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
F E D C B A Z Y X W V U T S R Q P O N M L K J I H G
F F

There are $3 \cdot 3 \cdot 4 \cdot 2 = 72$ pairs out of the most frequent characters.

Some of them are:

EIAAN ... ITAAN ... ITEIN ... ITINN ... ITISN
IEIAS ... ETIAS ... ETESS ... ETANS ... ETAIS

Defense against this attack: random key stream
 → one time pad

However, never use the same key twice. Otherwise:

$$(a_1, \dots, a_n) \oplus (k_1, \dots, k_n) = (c_1, \dots, c_n)$$

$$(b_1, \dots, b_n) \oplus (k_1, \dots, k_n) = (d_1, \dots, d_n)$$

Oscar:

$$(c_i - d_i) \bmod 26 = (a_i - b_i) \bmod 26$$

vulnerable to the above attack.

4. Entropy and Perfect Secrecy

4.1. Entropy

Consider random experiments, e.g.,

$$(0.9, 0.05, 0.05)$$

$$(0.33, 0.33, 0.34)$$

We aim at a measure of

= { uncertainty about the outcome (before)
information gained by the outcome (after)

The right measure was introduced by Shannon (1949).

Formal description

X : discrete r.v. with finite support $X = \{x_1, \dots, x_n\}$

distribution: $P(X=x_i) = p_i, i=1, \dots, n$

Def. 4.1. Let $c > 1$ constant.

$$H(X) = - \sum_{i=1}^n p_i \log_c p_i = - \sum_{i=1}^n P(X=x_i) \log_c P(X=x_i)$$

is called entropy of X (or (p_1, \dots, p_n)).

Convention: $0 \cdot \log 0 = 0$, omit c but fix it.

Analogously for 2-dim. random variables

(X, Y) with support $\mathcal{X} \times \mathcal{Y} = \{x_1, \dots, x_m\} \times \{y_1, \dots, y_d\}$
distribution $P(X=x_i, Y=y_j) = p_{ij}$

Def 4.2.

$$\begin{aligned} \text{a) } H(X, Y) &= - \sum_{i,j} P(X=x_i, Y=y_j) \log P(X=x_i, Y=y_j) \\ &= - \sum_{i,j} p_{ij} \log p_{ij} \end{aligned}$$

is called (joint) entropy of X, Y .

$$\begin{aligned} \text{b) } H(X|Y) &= - \sum_{j=1}^d P(Y=y_j) \sum_{i=1}^m P(X=x_i | Y=y_j) \log P(X=x_i | Y=y_j) \\ &= - \sum_{i,j} P(X=x_i, Y=y_j) \log P(X=x_i | Y=y_j) \end{aligned}$$

is called conditional entropy or equivocation. \perp

Theorem 4.3.

a) $0 \stackrel{(i)}{\leq} H(X) \stackrel{(ii)}{\leq} \log m$

"=" in (i) $\Leftrightarrow \exists x_i: P(X=x_i)=1$

"=" in (ii) $\Leftrightarrow P(X=x_i)=\frac{1}{m} \forall i$

b) $0 \stackrel{(i)}{\leq} H(X|Y) \stackrel{(ii)}{\leq} H(X)$

"=" in (i) $\Leftrightarrow P(X=x_i | Y=y_j)=1 \quad \forall i, j \text{ with } P(X=x_i, Y=y_j) > 0$

"=" in (ii) $\Leftrightarrow X, Y$ stoch. indep.

c) $H(X) \stackrel{(i)}{\leq} H(X, Y) \stackrel{(ii)}{\leq} H(X) + H(Y)$

"=" in (i) $\Leftrightarrow Y$ is fully dependent on X

"=" in (ii) $\Leftrightarrow X, Y$ stoch. independent

d) $H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$
(chain rule)

Proof. any book on information theory

↓