# 187 - Day 09

Day 09 – Breaking Vigenère cipher (cont.)

## Index of Coincidence

The following table give the relative frequency of the English alphabet letters in a 7834-letter sample of English writing.

| Letter | Relative frequency | Letter | Relative frequency |
|--------|--------------------|--------|--------------------|
| A | 0.08399 | N | 0.06778 |
| B | 0.01442 | O | 0.07493 |
| C | 0.02527 | P | 0.01991 |
| D | 0.04800 | Q | 0.00077 |
| E | 0.12150 | R | 0.06063 |
| F | 0.02132 | S | 0.06319 |
| G | 0.02323 | T | 0.08999 |
| H | 0.06025 | U | 0.02783 |
| I | 0.06485 | V | 0.00996 |
| J | 0.00102 | W | 0.02464 |
| K | 0.00689 | X | 0.00204 |
| L | 0.04008 | Y | 0.02157 |
| M | 0.02566 | Z | 0.00025 |

*[handwritten annotation:]* $\dfrac{\text{freq. of N}}{7834} = P_N$ — prob. that a randomly chosen letter in the plaintext is "N".

The probability that two randomly selected letters in English are identical is given by

$$\sum_{\alpha=A}^{Z} p_\alpha^2 \approx 0.065$$

In a Vigenère cipher with sufficiently long keyword, the probabilities of seeing any letter in the ciphertext will converge to

$$\frac{1}{26} = 0.0385$$

1

## Friedman Test

**Definition 1** (Index of Coincidence). *The **index of coincidence** (for a ciphertext), denoted $I$, is the probability that two randomly selected letters in the ciphertext are identical.*

Remark:

- If $I \approx 0.065$ then the cipher is more likely to be mono-alphabetic substitution.

- For poly-alphabetic substitution, $0.0385 \leq I \leq 0.065$

**Theorem 1.** *Let $n_0, n_1, n_2, \ldots n_{24}, n_{25}$ be the respective counts of the letters* (in ciphertext.) *$A, B, C, \ldots, Y, Z$. Let $n = \sum n_i$ be the total number of letters in the text then*

$$I = \frac{1}{n(n-1)} \sum_{i=0}^{25} n_i(n_i - 1).$$

*Now if an English plaintext is encrypted using a Vigenère cipher with keyword of length $k$, then*

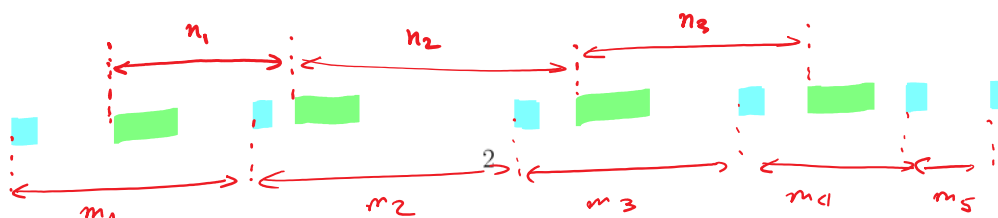$$I \approx \frac{0.0385 \cdot n(k-1) + 0.065(n-k)}{k(n-1)}, \text{ or equivalently,}$$

$(0.065 - 0.0385)$
↓
different for another language.

$$k \approx \frac{0.0265n}{(0.065 - I) + n(I - 0.0385)}.$$

## Kasiski Test

The **Kasiski Test** is another way of estimating the length of the keyword for Vigenère cipher. It obtains possible keyword lengths from the **gcd of the spacing between repeated letter groups** in the ciphertext.



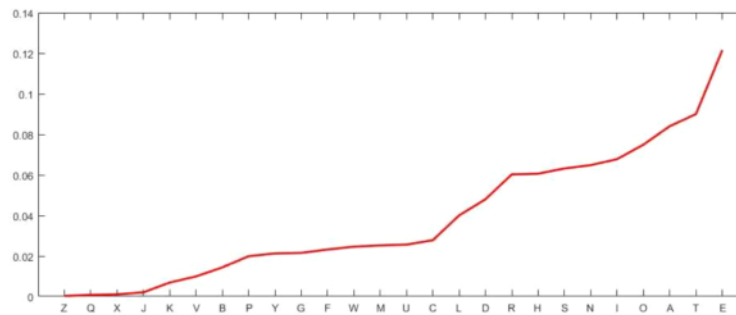$$k \approx \gcd\left(n_1, \ldots, n_3, m_1, m_2, \ldots, m_5\right)$$

# Cryptanalysis of Vigenère cipher

<u>Remark:</u> Both Friedman and Kasiski Tests only give the keyword length, but ) **bad.**
not the keyword itself. Firthermore, they are not very accurate when the )
ciphertext is small (usually less than 400 characters).

The **signature of English** is the graph of letter frequency distribution
of English when we sort these frequencies in increasing order.
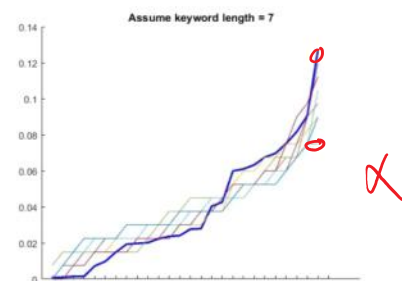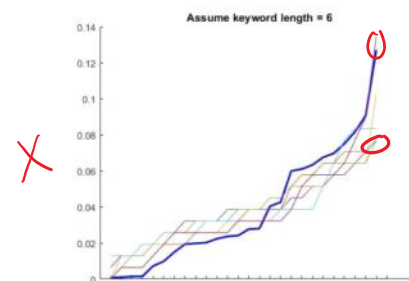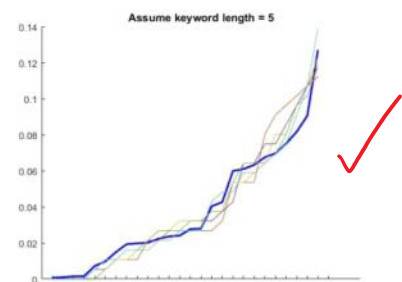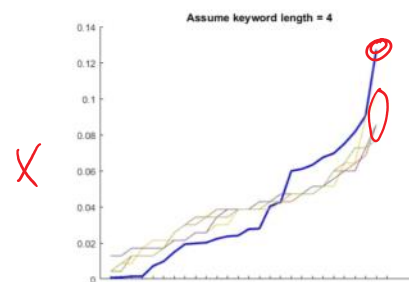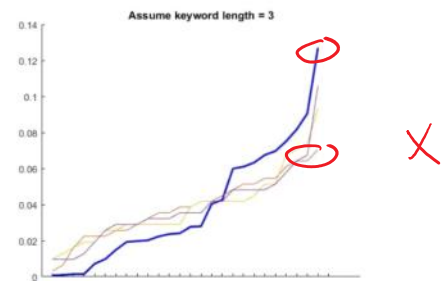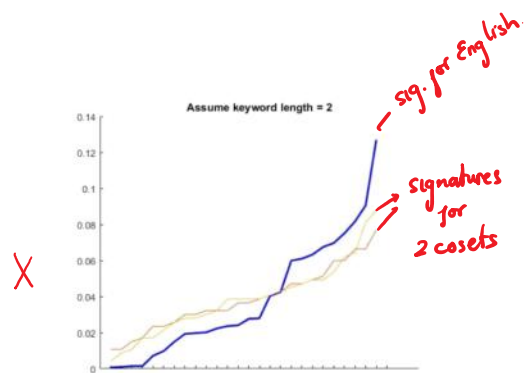
*To draw the signature :*

*· Obtain the frequency count for each letter*

*· sort the list in increasing order*

*· plot these points*

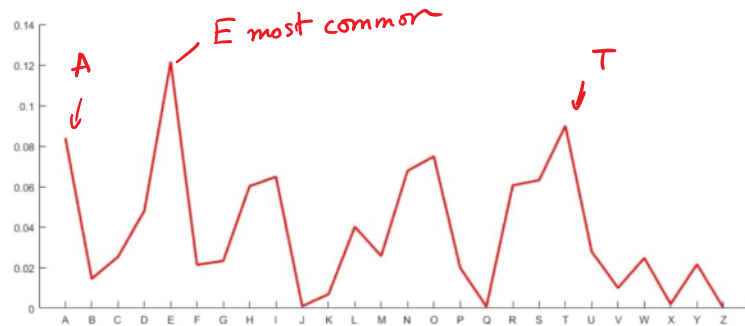

A **coset** are all the letters of the Vigenère ciphertext that are encrypted
by the same letter of the keyword.

<u>Remark:</u> A coset of the Vigenère ciphertext has the same encryption as
a Caesar shift cipher.

*· pick several values for keyword length k.*

*· draw the signature for all cosets*

*· If k is the correct value then the graphs will behave similar to the signature of English.*

Assume keyword length = 2

sig. for English.

signatures for 2 cosets

Assume keyword length = 3

Assume keyword length = 4

Assume keyword length = 5

Assume keyword length = 6

Assume keyword length = 7

4

The **scrawl of English** is the graph of letter frequency distribution of English in alphabetical order.



- Scrawl of each coset must be of the same shape. as above, except that the coset's scrawls may be shifted.

- To get the keyword: try to match the scrawl of each coset to that of English and record the shift distance

- look at the applet on Assignment page.

## Monty Hall Problem
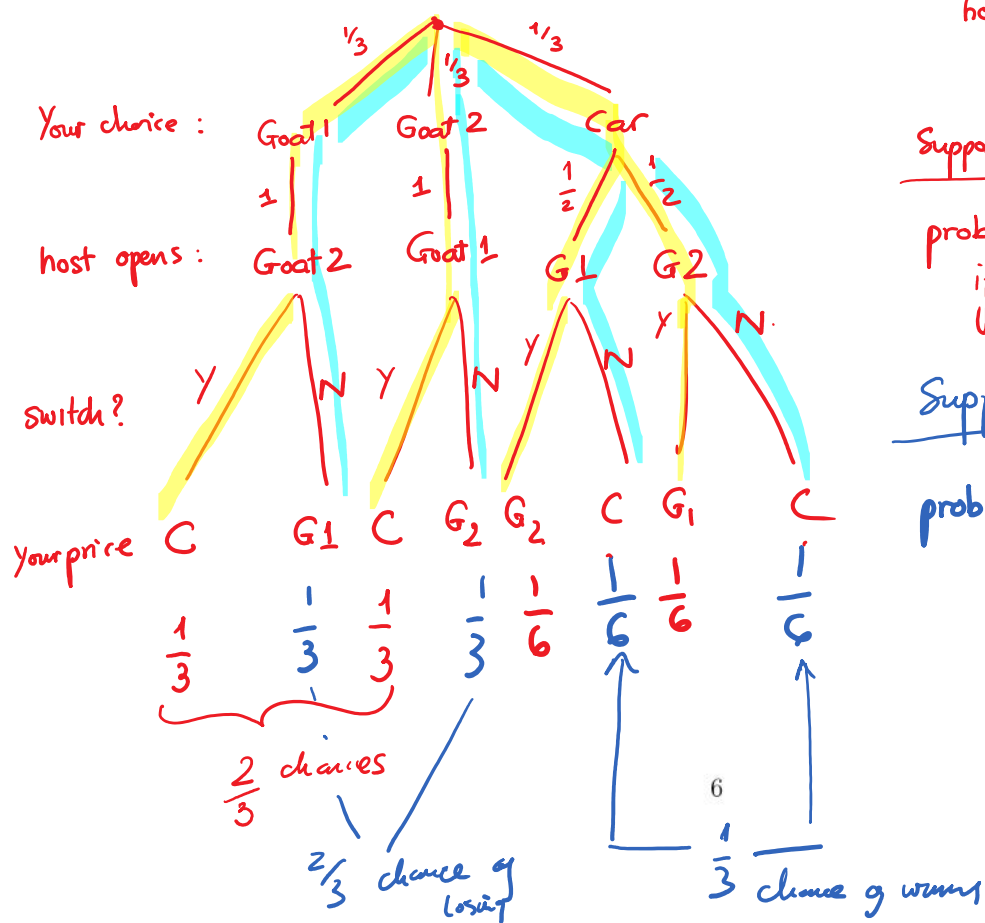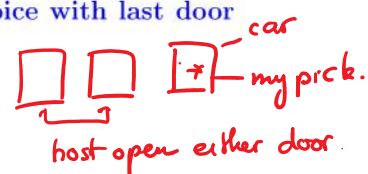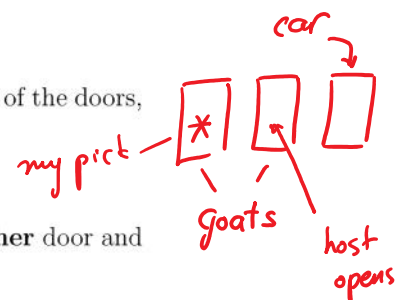
Consider the following game:

Suppose there are three doors. A car is hidden behind one of the doors, and the other two have goats.

As a player, you pick one of the doors.

The host, who **knows where the car is**, will open **another** door and reveal a goat.

Now you can choose whether to swap your choice with last door or stay with original choice

What is your best strategy?

car

my pick

goats

host opens

car

my pick.

host open either door.

Your choice :  Goat 1    Goat 2    Car

$\frac{1}{3}$    $\frac{1}{3}$    $\frac{1}{3}$

$1$    $1$    $\frac{1}{2}$  $\frac{1}{2}$

host opens :   Goat 2    Goat 1    G1    G2

switch?    Y    N    Y    N    Y    N    Y    N

Your prize   C    G1    C    G2    G2    C    G1    C

$\frac{1}{3}$    $\frac{1}{3}$    $\frac{1}{3}$    $\frac{1}{3}$    $\frac{1}{6}$    $\frac{1}{6}$    $\frac{1}{6}$    $\frac{1}{6}$

$\frac{2}{3}$ chances

$\frac{2}{3}$ chance of losing

6

$\frac{1}{3}$ chance of winning

Suppose you switch!

prob. that you win $= \frac{2}{3}$ if you switch

Suppose you stay:

prob. that you win $= \frac{1}{3}$ if you stay

**Definition 2.** *The **conditional probability** of an event B is the probability that this event will occur, given the knowledge that another event A has already occurred.*

$$\mathbb{P}(B|A) = \frac{\mathbb{P}(B \text{ and } A)}{\mathbb{P}(A)},$$

*assuming that* $\mathbb{P}(A) > 0$.

Refer to the table for the frequency of character pairs in English language.

Take an English text, pinpoint a letter at position $\lambda$

___ $\lambda$ ___

prob. that $\lambda$ is "A" ?

$$\mathbb{P}(\lambda = \text{"A"}) = p_A = 0.08399$$

Now if I know the letter $\mu$ to the left of $\lambda$ — .

how will the prob. change?

$$\overline{\mu \; \lambda}$$

What is the prob. $\lambda = \text{"A"}$, if $\mu$ is:

• $\mu = \text{"L"}$ ?

• $\mu = \text{"E"}$ ?

• $\mu = \text{"Q"}$ ?  } use conditional probability & today table.