

Lecture 2: Classical Encryption Techniques

Lecture Notes on “Computer and Network Security”

by Avi Kak (kak@purdue.edu)

January 14, 2010

©2010 Avinash Kak, Purdue University

Two Goals:

- To introduce the rudiments of encryption vocabulary.
- To trace the history of some early approaches to cryptography and to show through this history a common failing of humans to get carried away by the technological and scientific hubris of the moment.

2.1: Basic Vocabulary of Classical Encryption

plaintext: This is what you want to encrypt

ciphertext: The encrypted output

enciphering or encryption: The process by which plaintext is converted into ciphertext

encryption algorithm: The sequence of data processing steps that go into transforming plaintext into ciphertext. Various parameters used by an encryption algorithm are derived from a secret key.

In classical cryptography for commercial and other civilian applications, the encryption algorithm is made public.

secret key: A secret key is used to set some or all of the various parameters used by the encryption algorithm. **The important thing to note is that the same secret key is used for encryption and decryption in classical cryptography.** It is for this reason that classical cryptography is also referred to as **symmetric key cryptography.**

deciphering or decryption: Recovering plaintext from ciphertext

decryption algorithm: The sequence of data processing steps that go into transforming ciphertext back into plaintext. Various parameters used by a decryption algorithm are derived from the same secret key that was used in the encryption algorithm.

In classical cryptography for commercial and other civilian applications, the decryption algorithm is made public.

cryptography: The many schemes available today for encryption and decryption

cryptographic system: Any single scheme for encryption

cipher: A cipher means the same thing as a “cryptographic system”

block cipher: A block cipher processes a block of input data at a time and produces a ciphertext block of the same size.

stream cipher: A stream cipher encrypts data on the fly, usually one byte at a time.

cryptanalysis: Means “breaking the code”. Cryptanalysis relies on a knowledge of the encryption algorithm (that for civilian applications should be in the public domain) and some knowledge of the possible structure of the plaintext (such as the structure of a typical inter-bank financial transaction) for a partial or full reconstruction of the plaintext from ciphertext. Additionally, the goal is to also infer the key for decryption of future messages.

The precise methods used for cryptanalysis depend on whether the “attacker” has just a piece of ciphertext, or pairs of plaintext and ciphertext, how much structure is possessed by the plaintext, and how much of that structure is known to the attacker.

All forms of cryptanalysis for classical encryption exploit the fact that some aspect of the structure of plaintext may survive in the ciphertext.

brute-force attack: When encryption and decryption algorithms are publicly available, a brute-force attack means trying every possible key on a piece of ciphertext until an intelligible translation into plaintext is obtained.

key space: The total number of all possible keys that can be used in a cryptographic system. For example, **DES** uses a 56-bit key. So the key space is of size 2^{56} , which is approximately the same as 7.2×10^{16} .

cryptology: Cryptography and cryptanalysis together constitute the area of cryptology

2.2: Building Blocks of Classical Encryption Techniques

- Two building blocks of all classical encryption techniques are **substitution** and **transposition**.
- Substitution means replacing an element of the plaintext with an element of ciphertext.
- Transposition means rearranging the order of appearance of the elements of the plaintext.
- Transposition is also referred to as permutation.

2.3: Caesar Cipher

- This is the earliest known example of a substitution cipher.
- Each character of a message is replaced by a character three positions down in the alphabet.

plaintext: are you ready

ciphertext: DUH BRX UHGDB

- If we represent each letter of the alphabet by an integer that corresponds to its position in the alphabet, the formula for replacing each character 'p' of the plaintext with a character 'C' of the ciphertext can be expressed as

$$C = E(3, p) = (p + 3) \bmod 26$$

- A more general version of this cipher that allows for any degree of shift would be expressed by

$$C = E(k, p) = (p + k) \bmod 26$$

- The formula for decryption would be

$$p = D(k, C) = (C - k) \bmod 26$$

- In these formulas, 'k' would be the secret key. The symbols 'E' and 'D' represent encryption and decryption.

2.4: The Swahili angle ...

- A simple substitution cipher obviously looks much too simple, but that is the case only if you have some idea regarding the nature of the plaintext.
- What if the “plaintext” could be considered to be a binary stream of data and a substitution cipher replaced every consecutive 6 bits with one of 64 possible cipher characters? *In fact, this is referred to as Base64 encoding for sending email multimedia attachments.*
- If you did not know anything about the underlying plaintext and it was encrypted by a Base64 sort of algorithm, it might not be as trivial a cryptographic system as the substitution cipher shown on the previous page. But, of course, if the word ever got out that your plaintext was in Swahili, you’d be hosed.

2.5: Monoalphabetic Ciphers

- In a monoalphabetic cipher, our substitution characters are a random permutation of the 26 letters of the alphabet:

plaintext letters:	a	b	c	d	e	f
substitution letters:	t	h	i	j	a	b

- The key now is the sequence of substitution letters. In other words, the key in this case is the actual random permutation of the alphabet used.
- Note that there are $26!$ permutations of the alphabet. That is a number larger than 4×10^{26} .

2.6: A Very Large Key Space But

- That gives us a huge key space (meaning the total number of all possible keys that would need to be guessed in a brute-force attack). This key space is 10 orders of magnitude larger than the size of the key space for DES, the now somewhat outdated (but still widely used) NIST standard.
- Obviously, this would rule out a brute-force attack. (Even if each key took only a nanosecond to try, it would still take zillions of years to try out even half the keys.)
- So this would seem to be the answer to our prayers for an unbreakable code for symmetric encryption.
- But it is not!
- Why? Read on.

2.7: The All-Fearsome Statistical Attack

- If you know the nature of plaintext, any substitution cipher, regardless of the size of the key space, can be broken easily with a statistical attack.
- When the plaintext is plain English, a simple form of statistical attack consists measuring the frequency distribution for single characters, for pairs of characters, for triples of characters, etc., and comparing those with similar statistics for English.
- Figure 1 shows the relative frequency of of the letters in a sample of English text. Obviously, by comparing this distribution with a histogram for the characters in a piece of ciphertext, you may be able to establish the true identities of the ciphertext characters.

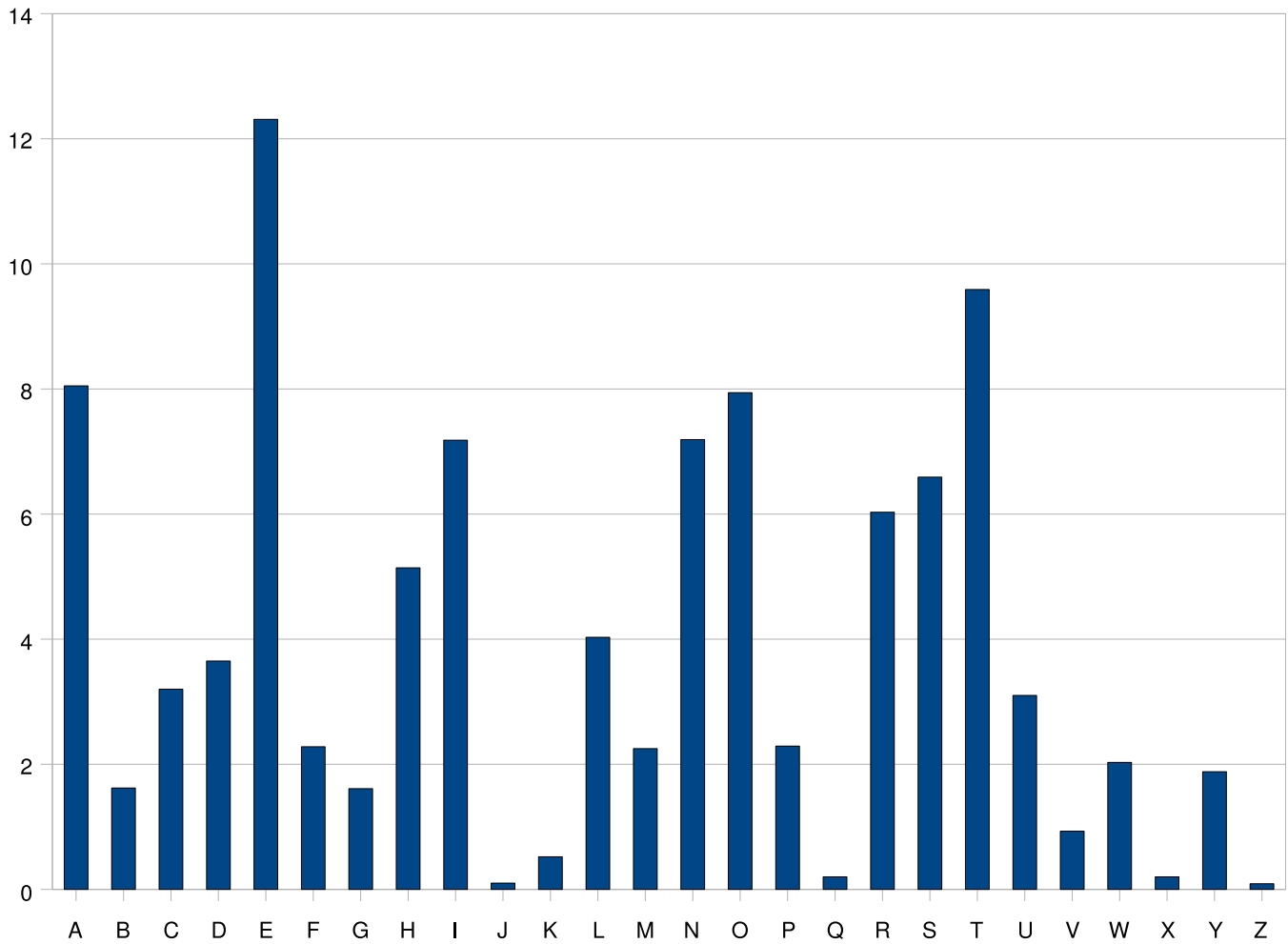


Figure 1: *This figure is from Lecture 2 of “Computer and Network Security” by Avi Kak*

2.8: Comparing the Statistics for Digrams and Trigrams

- Equally powerful statistical inferences can be made by comparing the relative frequencies for pairs and triples of characters in the ciphertext and the language believed to be used for the plaintext.
- Pairs of adjacent characters are referred to as **digrams**, and triples of characters as **trigrams**.
- Shown in Table 1 are the digram frequencies. The table does not include digrams whose relative frequencies are below 0.47. (A complete table of frequencies for all possible digrams would have 676 entries in it.)
- Let's say we have available to us the relative frequencies for all possible digrams. Let's represent this table by $p(x, y)$ where x denotes the first letter of a digram and y the second letter.

- The most frequently occurring trigrams ordered by decreasing frequency are:

the and ent ion tio for nde

<i>digram</i>	<i>frequency</i>	<i>digram</i>	<i>frequency</i>	<i>digram</i>	<i>frequency</i>	<i>digram</i>	<i>frequency</i>
th	3.15	to	1.11	sa	0.75	ma	0.56
he	2.51	nt	1.10	hi	0.72	ta	0.56
an	1.72	ed	1.07	le	0.72	ce	0.55
in	1.69	is	1.06	so	0.71	ic	0.55
er	1.54	ar	1.01	as	0.67	ll	0.55
re	1.48	ou	0.96	no	0.65	na	0.54
es	1.45	te	0.94	ne	0.64	ro	0.54
on	1.45	of	0.94	ec	0.64	ot	0.53
ea	1.31	it	0.88	io	0.63	tt	0.53
ti	1.28	ha	0.84	rt	0.63	ve	0.53
at	1.24	se	0.84	co	0.59	ns	0.51
st	1.21	et	0.80	be	0.58	ur	0.49
en	1.20	al	0.77	di	0.57	me	0.48
nd	1.18	ri	0.77	li	0.57	wh	0.48
or	1.13	ng	0.75	ra	0.57	ly	0.47

Table 1: *This table is from Lecture 2 of “Computer and Network Security” by Avi Kak*

2.9: Multiple-character Encryption to Mask Plaintext Structure

- One character at a time substitution obviously leaves too much of the plaintext structure in ciphertext.
- So how about destroying some of that structure by mapping multiple characters at a time to ciphertext characters?
- The best known approach that carries out multiple-character substitution is known as **Playfair cipher**.

2.10: Constructing the Matrix for Pairwise Substitutions in Playfair Cipher

In Playfair cipher, you first choose an encryption key. You then enter the letters of the key in the cells of a 5×5 matrix in a left to right fashion starting with the first cell at the top-left corner. You fill the rest of the cells of the matrix with the remaining letters in alphabetic order. The letters I and J are assigned the same cell. In the following example, the key is “**smythework**”:

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

2.11: Substitution Rules for Pairs of Characters in Playfair Cipher

1. Two plaintext letters that fall in the same row of the 5×5 matrix are replaced by letters to the right of each in the row. The “rightness” property is to be interpreted circularly in each row, meaning that the first entry in each row is to the right of the last entry. Therefore, the pair of letters “bf” in plaintext will get replaced by “CA” in ciphertext.
2. Two plaintext letters that fall in the same column are replaced by the letters just below them in the column. The “belowness” property is to be considered circular, in the sense that the topmost entry in a column is below the bottom-most entry. Therefore, the pair “ol” of plaintext will get replaced by “CV” in ciphertext.
3. Otherwise, for each plaintext letter in a pair, replace it with the letter that is in the same row but in the column of the other letter. Consider the pair “gf” of the plaintext. We have ‘g’ in the fourth row and the first column; and ‘f’ in the third row and the fifth column. So we replace ‘g’ by the letter in the same row as ‘g’ but in the column that contains ‘f’. This gives us ‘P’ as a replacement for ‘g’. And we replace ‘f’ by the letter in the same row as ‘f’ but in the column that contains ‘g’. That gives us ‘A’

as replacement for 'f'. Therefore, 'gf' gets replaced by 'PA'.

2.12: Dealing with Duplicates Letters in a Key and Repeating Letters in Plaintext

- You must drop any duplicates in a key.
- Before the substitution rules are applied, you must insert a chosen “filler” letter (let’s say it is ‘x’) between any repeating letters in the plaintext. So a plaintext word such as “hurray” becomes “hurxray”

2.13: How Secure is the Playfair Cipher?

- Playfair was thought to be unbreakable for many decades.
- It was used as the encryption system by the British Army in World War 1. It was also used by the U.S. Army and other Allied forces in World War 2.
- But, as it turned out, Playfair was extremely easy to break.
- As expected, the cipher does alter the relative frequencies associated with the individual letters and with digrams and with trigrams, but not sufficiently.
- The figure on page 22 shows the single-letter relative frequencies in descending order (and normalized to the relative frequency of the letter 'e') for different ciphers. There is still considerable information left in the distribution for good guesses.
- The cryptanalysis of the Playfair cipher is also aided by the fact that a digram and its reverse will encrypt in a similar fashion.

That is, if AB encrypts to XY, then BA will encrypt to YX. So by looking for words that begin and end in reversed digrams, one can try to compare them with plaintext words that are similar. Example of words that begin and end in reversed digrams: receiver, departed, repairer, redder, denuded, etc.

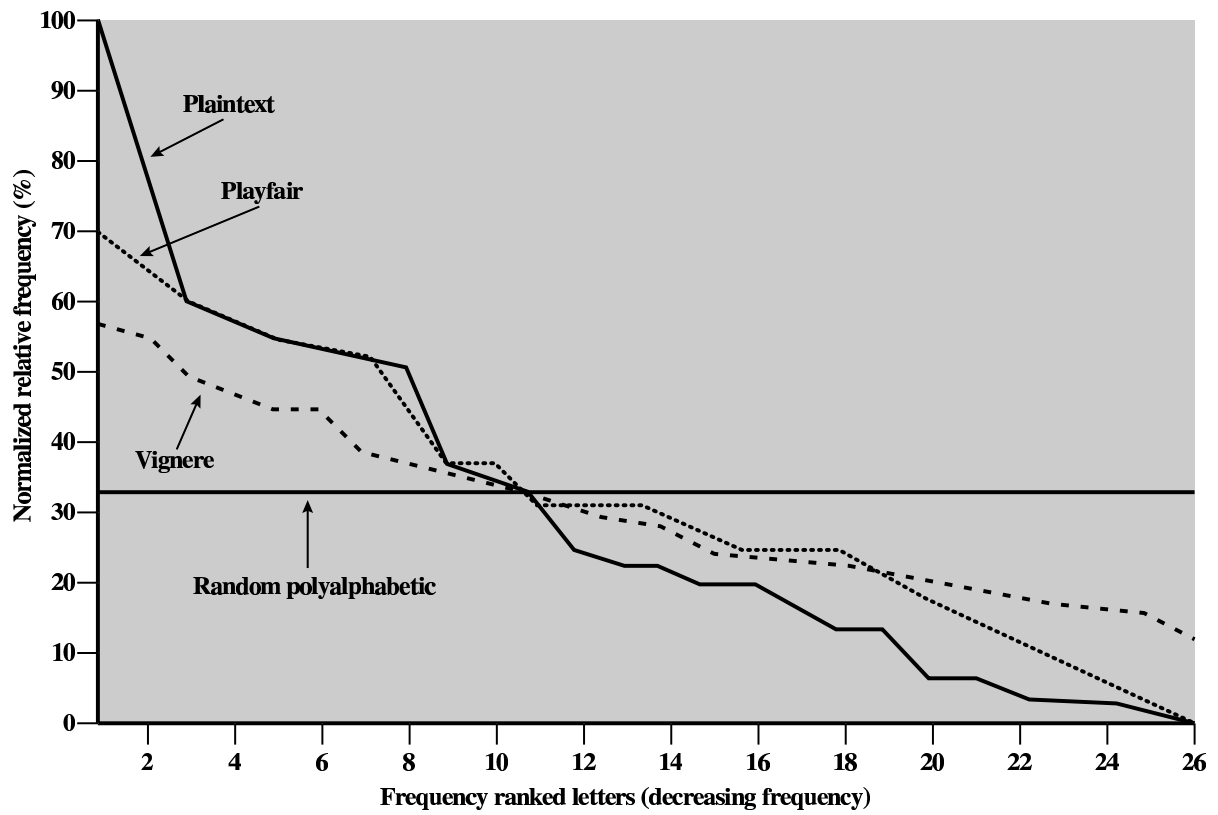


Figure 2.6 Relative Frequency of Occurrence of Letters

This figure is from Chapter 2 of William Stallings: “Cryptography and Network Security”, Fourth Edition, Prentice-Hall.

2.14: Another Multi-Letter Cipher: The Hill Cipher

The Hill cipher takes a very different (more mathematical) approach to multi-letter substitution:

- You assign an integer to each letter of the alphabet. For the sake of discussion, let's say that you have assigned the integers 0 through 25 to the letters 'a' through 'z' of the plaintext.
- The encryption key, call it \mathbf{K} , consists of a 3×3 matrix of integers:

$$\mathbf{K} = \begin{bmatrix} k_{11} & k_{12} & k_{13} \\ k_{21} & k_{22} & k_{23} \\ k_{31} & k_{32} & k_{33} \end{bmatrix}$$

- Now we can transform three letters at a time from plaintext, the letters being represented by the numbers p_1 , p_2 , and p_3 , into three ciphertext letters c_1 , c_2 , and c_3 in their numerical representations by

$$\begin{aligned} c_1 &= (k_{11}p_1 + k_{12}p_2 + k_{13}p_3) \text{ mod } 26 \\ c_2 &= (k_{21}p_1 + k_{22}p_2 + k_{23}p_3) \text{ mod } 26 \\ c_3 &= (k_{31}p_1 + k_{32}p_2 + k_{33}p_3) \text{ mod } 26 \end{aligned}$$

- The above set of linear equations can be written more compactly in the following vector-matrix form:

$$\vec{\mathbf{C}} = [\mathbf{K}] \vec{\mathbf{P}} \text{ mod } 26$$

- Obviously, the decryption would require the inverse of \mathbf{K} matrix.

$$\vec{\mathbf{P}} = [\mathbf{K}^{-1}] \vec{\mathbf{C}} \text{ mod } 26$$

This works because

$$\vec{\mathbf{P}} = [\mathbf{K}^{-1}] [\mathbf{K}] \vec{\mathbf{P}} \text{ mod } 26 = \vec{\mathbf{P}}$$

2.15: How Secure is the Hill Cipher?

- It is extremely secure against ciphertext only attacks. That is because the keyspace can be made extremely large by choosing the matrix elements from a large set of integers. (The key space can be made even larger by generalizing the technique to larger-sized matrices.)
- But it has zero security when the plaintext–ciphertext pairs are known. The key matrix can be calculated easily from a set of known $\vec{\mathbf{P}}, \vec{\mathbf{C}}$ pairs.

2.16: Polyalphabetic Ciphers: The Vigenere Cipher

- In a monoalphabetic cipher, the same substitution rule is used for every substitution. In a polyalphabetic cipher, the substitution rule changes continuously from letter to letter according to the elements of the encryption key.
- Let each letter of the encryption key denote a shifted Caesar cipher, the shift corresponding to the key. This is illustrated with the help of the table on the next page.
- Now a plaintext message may be encrypted as follows

key:	abracadabraabracadabraabracadabraab
plaintext:	canyoumeetmeatmidnightihavethegoods
ciphertext:	CBEYQUPEFKMEBK.....

- The Vigenere cipher is an example of a polyalphabetic cipher.
- Since, in general, the encryption key will be shorter than the message to be encrypted, for the Vigenere cipher the key is repeated, as illustrated in the above example where the key is the string “abracadabra”.

<i>encryption key</i>	<i>plain text letters</i>				
<i>letter</i>	a	b	c	d
	<i>substitution letters</i>				
<i>a</i>	A	B	C	D
<i>b</i>	B	C	D	E
<i>c</i>	C	D	E	F
<i>d</i>	D	E	F	G
<i>e</i>	E	F	G	H
.
.
<i>z</i>	Z	A	B	C

2.17: How Secure is the Vigenere Cipher?

- Since there exist in the output multiple ciphertext letters for each plaintext letter, you would expect that the relative frequency distribution would be effectively destroyed. But as can be seen in the plots on page 22, a great deal of the input statistical distribution still shows up in the output. (The plot shown for Vigenere cipher is for an encryption key that is 9 letters long.)
- Obviously, the longer the encryption key, the greater the masking of the structure of the plaintext. The best possible key is as long as the plaintext message and consists of a purely random permutation of the 26 letters of the alphabet. This would yield the ideal plot shown in the figure on page 22 of these notes. The ideal plot is labeled “Random polyalphabetic” in the figure.
- In general, to break the Vigenere cipher, you first try to estimate the length of the encryption key. This length can be estimated by using the logic that plaintext words separated by multiples of the length of the key will get encoded in the same way.
- If the estimated length of the key is N , then the cipher consists of N monoalphabetic substitution ciphers and the plaintext letters at positions 1, N , $2N$, $3N$, etc., will be encoded by the same

monoalphabetic cipher. This insight can be useful in the decoding of the monoalphabetic ciphers involved.

2.18: Transposition Techniques

- All of our discussion so far has dealt with substitution ciphers. We have talked about monoalphabetic substitutions, polyalphabetic substitutions, etc.
- We will now talk about a different notion in classical cryptography: permuting the plaintext.
- This is how a pure permutation cipher could work: You write your plaintext message along the rows of a matrix of some size. You generate ciphertext by reading along the columns. The order in which you read the columns is determined by the encryption key:

key: 4 1 3 6 2 5

plaintext: m e e t m e
 a t m i d n
 i g h t f o
 r t h e g o
 d i e s x y

ciphertext: ETGTIMDFGXEMHHEMAIRDENOOYTITES

- The cipher can be made more secure by performing multiple rounds of such permutations.

HOMEWORK PROBLEMS

1. All classical ciphers are based on symmetric key encryption. What does that mean?
2. What are the two building blocks of all classical ciphers?
3. True or false: The larger the size of the key space, the more secure a cipher? Justify your answer.
4. Give an example of a cipher that has an extremely large key space size, an extremely simple encryption algorithm, and extremely poor security.
5. What is the difference between monoalphabetic substitution ciphers and polyalphabetic substitution ciphers?
6. What is the main security flaw in the Hill cipher?
7. What makes Vigenere cipher more secure than, say, the Playfair cipher?
8. **Programming Assignment:**
Write a script called `hist.pl` in Perl (or `hist.py` in Python) that makes a histogram of the letter frequencies in a text file. The output should look like

```
A: xx
B: xx
C: xx
...
...
```

where xx stands for the count for that letter.

9. Programming Assignment:

Write a script called `poly_cipher.pl` in Perl (or `poly_cipher.py` in Python) that is an implementation of the Vigenere polyalphabetic cipher for messages composed from the letters of the English alphabet, the numerals 0 through 9, and the punctuation marks ‘.’, ‘,’, and ‘?’.

Your script should read from standard input and write to standard output. It should prompt the user for the encryption key.

Your hardcopy submission for this homework should include some sample plaintext, the ciphertext, and the encryption key used.

Make your scripts as compact and as efficient as possible. Make liberal use of builtin functions for what needs to be done. For example, you could make a circular list with either of the following two constructs in Perl:

```
unshift( @array, pop(@array) )
push( @array, shift(@array) )
```

See `perlfaq4` for some tips on array processing in Perl.

CREDITS

The data presented in Figure 1 and Table 1 are from <http://jnicholl.org/Cryptanalysis/Data/EnglishData.php>. That site also shows a complete digram table for all 676 pairings of the letters of the English alphabet.