

Encryption assignment

Izaak van Dongen

January 18, 2018

Contents

1	The Vigenère cipher (UEncrypt.pas)	1
2	Reading files (UFiles.pas)	2
3	Command line interface (PVigenere.pas)	3
4	Compiling and testing	3
5	Determining keyword length	5
6	Determining the keyword	7
7	Cracking a cipher	8
8	Source	9

The Vigenère cipher (UEncrypt.pas)

This section will concern the unit UEncrypt.pas, which contains my utility functions for Vigenère encryption. The unit's interface looks like this:

```
1 unit UEncrypt;
2
3 interface
4
5 uses SysUtils;
6
7 const
8     alph_s: string = 'ABCDEFGHJKLMNOPQRSTUVWXYZ';
9     alpha: set of char = ['A'.. 'Z', 'a'.. 'z'];
10    upper: set of char = ['A'.. 'Z'];
11
12 function vigenere(pt: ansistring; pass: string): ansistring;
13 function un_vigenere(pt: ansistring; pass: string): ansistring;
14 function proper_mod(a, b: integer): integer;
```

Listing 1: UEncrypt interface

The Vigenère cipher is a generalisation of the Caesar cipher. Because of this, in fact if I implement a Vigenère routine, I will also have implemented the Caesar shift algorithm. It would also implement the one-time pad, as this is simply a Vigenère cipher with a plaintext-length key. I will write my code to work on strings, as strings are guaranteed to be more easily 'seekable'. This is important as for later programs, I'm planning to do a lot of backwards-and-forwards reading through the ciphertext. Ideally I would have written the routine to work with something like 'an iterable of characters', but I have to make some compromises in Pascal.

A short-lived byproduct of this was that any input files were mysteriously cut off. Inspection with wc showed that each file was cut off to exactly 255 characters. Apparently, in some fpc modes 'string' is aliased to 'shortstring', which can only store so many characters. Because of this, I explicitly use 'ansistring' to represent plaintexts and ciphertexts, for clarity, rather than having to use specific compiler options.

Because the Vigenère algorithm is so fundamentally modular, I also needed to write a proper modulo function, such as in Python (that guarantees a positive result for any input). It looks like this:

```
1 function proper_mod(a, b: integer): integer;
2 begin
3     proper_mod := a mod b;
4     if proper_mod < 0 then
5         proper_mod := proper_mod + b;
```

```
6 end;
```

Listing 2: Proper modulo function

Because this function complies with the strict (useful) definition of modulo, this means the later encryption routines that rely on it can be used with any combination of characters, not just words, without worrying about going out of bounds. Note that almost all of them ‘rely’ on it as the following function will use it.

I also wrote a function to retrieve the ‘alphabetic ordinal’ of a letter (ie an integer $k:0 \leq k < 26$):

```
1 function alpha_ord(c: char): integer;
2 begin
3   alpha_ord := proper_mod(ord(upcase(c)) - 65, 26);
4 end;
```

Listing 3: Proper modulo function

Note that the previous two functions both behave in a ‘0-indexed’ kind of way - this is because the Vigenère cipher relies heavily on modulo and 0-indexing (A is considered to have the value 0, for example, and the wrapping is implemented by modulo). Because of this, when I’m working with Pascal’s actual (1-indexed) strings, careful consideration is needed to avoid off-by-one errors.

Anyway, here’s the moment you’ve all been waiting for - the Vigenère routine:

```
1 function vigenere(pt: ansistring; pass: string): ansistring;
2 var
3   i: integer = 0;
4   c: char;
5 begin
6   vigenere := '';
7   for c in pt do
8     if c in alpha then begin
9       if c in upper then
10        vigenere := vigenere + chr(65 + (alpha_ord(c)
11          + alpha_ord(pass[i + 1])) mod 26)
12      else
13        vigenere := vigenere + chr(97 + (alpha_ord(c)
14          + alpha_ord(pass[i + 1])) mod 26);
15      i := (i + 1) mod length(pass);
16    end else
17      vigenere := vigenere + c;
18 end;
```

Listing 4: Vigenère algorithm

It works quite simply by iterating over the plaintext and whenever it meets a character, advancing the position in the key, and applying the shift using some modular arithmetic, using a conditional to preserve case, and correcting for the flaw that is 1-indexing.

This function uses its output slot to dynamically accumulate the output, rather than making another variable of the same type as the return type named ‘answer’ or ‘ct’ or something. This approach will be used by several other functions.

Now, having written a Vigenère routine, it can in fact be recycled to write a decryption routine, by inverting the key, followed by delegation to the encryption routine. In this case, ‘inverting the key’ means finding the modular additive inverse of each component.

```
1 function un_vigenere(pt: ansistring; pass: string): ansistring;
2 var
3   i: integer;
4 begin
5   for i := 1 to length(pass) do
6     pass[i] := alph_s[1 + proper_mod(-alpha_ord(pass[i]), 26)];
7   un_vigenere := vigenere(pt, pass);
8 end;
```

Listing 5: Vigenère decryption

Reading files (UFiles.pas)

Having done the more fun part, we now have to briefly divert our attention to file handling. This unit handles all of the file reading, providing an interface to a single simple function:

```
1 function read_file(var f: textfile): ansistring;
2 var
3   c: char;
4 begin
5   read_file := '';
```

```

6   while not eof(f) do begin
7       read(f, c);
8       read_file := read_file + c;
9   end;
10 end;

```

Listing 6: File-reading routine

Command line interface (PVigenere.pas)

At last! A program that will do something. PVigenere will use the previous routines and interface with the command line to allow them to be used. The program is as follows:

```

1  {$MODE OBJFPC}
2
3  program PVigenere;
4
5  uses
6      UEncrypt, UFiles, SysUtils;
7
8  type
9      crypt_func = function(pt: ansistring; pass: string): ansistring;
10
11  var
12      pass: string;
13      arg_i: integer;
14      encryptor: crypt_func = @vigenere;
15  begin
16      pass := 'n';
17      for arg_i := 1 to paramcount do
18          if paramstr(arg_i) = '--decrypt' then
19              encryptor := @un_vigenere
20          else
21              try
22                  pass := alph_s[1 + proper_mod(strtoint(paramstr(arg_i)), 26)];
23              except
24                  on EConvertError do
25                      pass := paramstr(arg_i);
26              end;
27      write(output, encryptor(read_file(input), pass));
28  end.

```

Listing 7: Command-line interface for Vigenère routines (PVigenere.pas)

This ‘parses’ the command-line arguments to determine what to do with what parameters. If it encounters the argument ‘--decrypt’ anywhere, it uses the Vigenère decryption routine rather than the encryption routine, which it does by using a function variable to store the encryption routine, as both the encryptor and decryptor have the same signature.

The way it determines the key is as follows: first, the key is set to the default value of ‘n’. This results in a Caesar-shift of 13, or a ROT13 cipher. Then, for each argument that is not ‘--decrypt’, it sets to key to this argument, in some sense. First, it tries to parse the argument as an integer caesar shift. If this fails, the argument is used as a Vigenère key. For anything other than the alphabet this behaviour may not be very obvious, as it uses ASCII - 65 mod 26. However, the function is still well-defined and if the same key is used for decryption, it will successfully decrypt.

Because it must ‘attempt’ to parse an integer, it uses a try/catch statement. Because of this, the mode OBJFPC must be used.

Note also that this program is configured to work with an output file - although that output file is STDOUT, it still writes to a file stream rather than using the primitive write or writeln with no file parameter. Input is similarly read from STDIN. This is pretty standard for a command line program - if the user wants to read from actual input files, this can be easily done with `$ bin/Vigenere key < input.txt > output.txt`, or alternatively they might use the cat command. Leaving to writing to STDOUT is a lot easier for dynamic testing purposes, and doesn’t sacrifice any functionality.

Compiling and testing

This time I wrote another makefile, but one that works generically with my Pascal naming conventions. It assumes that any some `P(.*)\.pas` will compile to `bin/\1`, and also has a dependency on all local units, which take the form `U.*\.pas`.

```

1  .PHONY: all
2  all: $(shell echo P*.pas | sed -E "s/P([^.]*)\.pas/bin\/\1/g")
3
4  bin/%: P%.pas $(wildcard U*.pas)

```

```
5 fpc -TLinux -o$@ $<
```

Listing 8: The generic FPC makefile

Note that this in fact uses three different languages to represent a generic file - firstly, the generic makefile `bin/%` and `P%.pas`, secondly the shell globs `U*.pas` and `P*.pas`, and lastly the extended regex sed command `s/P([^.]*)\.pas/bin\/\1/g`. The make generics are needed to represent the instructions to build any specific file. The shell glob is used to find all possible program files in a directory, and the regex is used to scrape the executable names from the program file names. Now, I can run ‘make’, which builds the ‘all’ target, which is phony so will always try to build all dependencies, which in this case is just ‘bin/Vigenere’.

Anyway, having built my executable ‘bin/Vigenere’, I could now verify that it worked correctly. The first test string I used was `AbCdEfGhIjKlMnOpQrStUvWxYz`, as this would be useful to test if shifts were working correctly, and if case was being preserved. Don’t worry, I didn’t type it out by hand.

```
1 In [2]: """ .join(i.lower() if ind & 1 else i for ind, i in enumerate(string.ascii_uppercase))
2 Out [2]: 'AbCdEfGhIjKlMnOpQrStUvWxYz'
```

Listing 9: Alphabet

Anyway,

```
1 $ echo "AbCdEfGhIjKlMnOpQrStUvWxYz" | bin/Vigenere 1
2 BcDeFgHiJkLmNoPqRsTuVwXyZa
3 $ echo "AbCdEfGhIjKlMnOpQrStUvWxYz" | bin/Vigenere ab
4 AcCeEgGilkKmMoOqQsSuUwWyYa
5 $ echo "AbCdEfGhIjKlMnOpQrStUvWxYz" | bin/Vigenere 3
6 DeFgHiJkLmNoPqRsTuVwXyZaBc
7 $ echo "AbCdEfGhIjKlMnOpQrStUvWxYz" | bin/Vigenere 1 | bin/Vigenere --decrypt 1
8 AbCdEfGhIjKlMnOpQrStUvWxYz
9 $ echo "AbCdEfGhIjKlMnOpQrStUvWxYz" | bin/Vigenere ab | bin/Vigenere --decrypt ab
10 AbCdEfGhIjKlMnOpQrStUvWxYz
```

Listing 10: Testing bin/Vigenere

So far, so good. To really put it to the test, I decided to try it with the works of Shakespeare, on the assumption that all the edge cases would probably be covered somewhere along the way.

```
1 $ cat ~/programmeren/A453/text/shakespeare.txt | wc -c
2 5458199
3 $ time cat ~/programmeren/A453/text/shakespeare.txt | bin/Vigenere "William Shakespeare" | wc -c
4 5458199
5 cat ~/programmeren/A453/text/shakespeare.txt 0.00s user 0.02s system 5% cpu 0.399 total
6 bin/Vigenere "William Shakespeare" 0.84s user 0.07s system 99% cpu 0.920 total
7 wc -c 0.00s user 0.01s system 1% cpu 0.919 total
8 $ time diff <(cat ~/programmeren/A453/text/shakespeare.txt |
9 > bin/Vigenere "William Shakespeare" |
10 > bin/Vigenere --decrypt "William Shakespeare")
11 > ~/programmeren/A453/text/shakespeare.txt
12 diff ~/programmeren/A453/text/shakespeare.txt 0.00s user 0.02s system 1% cpu 1.792 total
```

Listing 11: Shakespeare

This actually did lead to the discovery a bug: as the works of Shakespeare are longer than the size of an integer, the integer (i) tracking the index in letters would overflow, leading to misalignment of the key and the plaintext. This was fixed by applying a modulus to the index whenever modifying it, rather than only when using it to index. See code listing 4. This has now been fixed, and as you can see, the ‘diff’ command with the source text and the decrypted plaintext exits cleanly, indicating it has been perfectly replicated.

As this is simultaneously an implementation of the one-time pad, we can use it as such. For now it will have to be more of a proof of concept, as the key is limited to <255 characters, as it’s stored in a string and taken from argv.

```
1 $ echo "the quick brown fox jumps over the lazy dog" | bin/Vigenere $(cat *.())
2 xuh xyvfr ferdr sre nhpww bylv gkl pncf hbj
3 $ echo "the quick brown fox jumps over the lazy dog" | bin/Vigenere $(cat *.()) | bin/Vigenere --decrypt
4 $(cat *.())
5 the quick brown fox jumps over the lazy dog
6 $ echo "the quick brown fox jumps over the lazy dog" | bin/Vigenere $(cat /dev/urandom | head -c 100)
7 jih shkno ugein gyr pwsce kgiu ycu mdbl fzk
```

Listing 12: Using the Vigenère program for a one-time pad

Here I’ve first used the concatenation of all the files in the directory as a recoverable key, to demonstrate also the decryption, followed by 100 cryptographically secure random bytes from `/dev/urandom` resulting in an unbreakable (undecryptable) encryption.

Determining keyword length

How, having implemented the easy approach for decryption (where you know the keyword), I decided to write a program that, for the user's convenience, doesn't require them to remember their Vigenère key. The first step is to determine the length of the keyword. This is absolutely trivial using some basic statistics and a tiny amount of computing power. For each hypothesised key length, we can simply consider each sequence of letters that would be encrypted by the same keyword letter - ie for a keyword of length l we get l different sub-texts of the ciphertext, where the i th of these is given by $S_i = \{C_{nl+i} : n \in \mathbb{N} \wedge nl+i < |c|\}$. If the hypothesis is correct, these will each have been individually encrypted with a single letter of the key, so we then calculate the index of coincidence of the distribution of each of these sub-texts. The IOC is the probability of any two letters being the same. This has the very obvious property that if the letters change around, the IOC will not change, so the IOC of a text is invariant under any kind of monoalphabetic substitution cipher. It can be calculated from a distribution d indexed from 1 to 26 as follows:

$$\text{IOC} = \frac{\sum_{i=1}^{26} d_i(d_i-1)}{T(T-1)}$$

where $T = \sum_{i=1}^{26} d_i$

The IOC is a very strong indicator of natural language/English. The expected IOC of English is around 0.067, whereas for uniformly distributed text the expected distribution is $\frac{1}{26} = 0.0385$.

All of the cracking functions have been written in one unit, 'UAttack.pas'. This unit has the following interface:

```
1 unit UAttack;
2
3 interface
4
5 uses UFiles , Math;
6
7 type
8     norm_dist = array ['A'..'Z'] of real;
9     dist = array ['A'..'Z'] of integer;
10
11 const
12     eng_dist: norm_dist =
13     (0.08167, 0.01492, 0.02782, 0.04253, 0.12702, 0.02228, 0.02015, 0.06094,
14     0.06966, 0.00153, 0.00772, 0.04025, 0.02406, 0.06749, 0.07507, 0.01929,
15     0.00095, 0.05987, 0.06327, 0.09056, 0.02758, 0.00978, 0.02360, 0.00150,
16     0.01974, 0.00074);
17     alpha: set of char = ['A'..'Z', 'a'..'z'];
18
19 function fitness(pt_dist: norm_dist): real;
20 function IOC(d: dist): real;
21 function get_dist(pt: ansistring): dist;
22 function get_interval(pt: ansistring; start, interval: integer): ansistring;
23 function clean(pt: ansistring): ansistring;
24 function normalise(d: dist): norm_dist;
```

Listing 13: UAttack interface

All of these functions will be covered in due course. Things to note are the 'dist' and 'norm_dist' types, which represent a discrete distribution and a normalised probability distribution respectively (ie 'dist' directly represents letter frequencies whereas 'norm_dist' is normalised so that its sum is 1). 'eng_dist' is the distribution of English, taken from Wikipedia.

Now, the first thing you need to do if you want to perform this analysis by IOC is to actually get our S_i substrings. For this, first we'll want to strip away all characters we aren't interested in. This is what the 'clean' function does:

```
1 function clean(pt: ansistring): ansistring;
2 var
3     c: char;
4 begin
5     clean := '';
6     for c in pt do
7         if c in alpha then
8             clean := clean + c;
9 end;
```

Listing 14: Clean function

Now, we can extract the sub-text, calculate its distribution, and calculate its distribution's IOC. This part is where the use of strings rather than file streams becomes really crucial, as we want to separately examine different parts of the strings, repeatedly.

```

1 function get_interval(pt: ansistring; start, interval: integer): ansistring;
2 var
3     i: integer;
4 begin
5     get_interval := '';
6     for i := 0 to (length(pt) - start) div interval do
7         get_interval := get_interval + pt[i * interval + start + 1];
8 end;
9
10 function get_dist(pt: ansistring): dist;
11 var
12     c: char;
13 begin
14     for c := 'A' to 'Z' do
15         get_dist[c] := 0;
16     for c in pt do
17         if c in alpha then
18             inc(get_dist[uppercase(c)]);
19 end;
20
21 function IOC(d: dist): real;
22 var
23     total: integer = 0;
24     freq: integer;
25 begin
26     IOC := 0;
27     for freq in d do begin
28         IOC := IOC + freq * (freq - 1);
29         total := total + freq;
30     end;
31     IOC := IOC / ((total * (total - 1)));
32 end;

```

Listing 15: The rest of the owl (IOC functions)

The only thing of real note here is that the IOC function only makes one pass, where it accumulates its denominator and the total T value.

Now, I needed to write a program that put these functions to use.

```

1 program PIOC;
2
3 uses UAttack, UFiles, Sysutils, Strutils;
4
5 function interval_ioc(pt: ansistring; interval: integer): real;
6 var
7     start: integer;
8 begin
9     interval_ioc := 0;
10    for start := 0 to interval - 1 do
11        interval_ioc := interval_ioc + IOC(get_dist(get_interval(pt, start, interval)));
12    interval_ioc := interval_ioc / interval;
13 end;
14
15 procedure print_IOC(pt: ansistring; max_interval: integer);
16 var
17     interval: integer;
18 begin
19     for interval := 1 to max_interval do begin
20         writeln(format('%2d (%.4f) %s', [interval,
21             interval_ioc(pt, interval),
22             dupestring('-',
23                 trunc(interval_ioc(pt, interval) * 500))]]);
24     end;
25 end;
26
27 begin
28     print_ioc(clean(read_file(input)), 20);
29 end.

```

Listing 16: Command-line interface to IOC functions (PIOC.pas)

This firstly calculates the average IOC for each of the l sub-texts for a hypothesised key length l , and then displays this as a bar

graph of l against IOC. My actual source code uses the character U+2796, or ‘HEAVY MINUS SIGN’, to produce a smooth bar graph. This uses the C-like ‘format’ function to ensure that everything is properly aligned (see listing 19).

Determining the keyword

Now that we can assume we know the length of the keyword, we can start to attack the keyword itself. To do this, we will need a little brute force - with very low complexity. The number of possibilities that need to be checked as 26^l , which is effectively constant, as for most ciphertexts l will be < 20 . Therefore the only real factor will be the length of the ciphertext, in which this algorithm should be roughly linear.

To fully automate this, we will need to produce a fitness metric for some text. A very good metric for this is quadgram probability - for each quadgram in the text, use a lookup table of probabilities for quadgrams to occur in English, and use these to accumulate a score for the text. I’ve used this previously with simulated annealing to attack ciphers that take permutations of the alphabet as a key (eg substitution, playfair, bifid). However, this is not suitable for this algorithm because this isn’t a holistic attack - we are attacking separate subsets of the text. Because of this, we won’t actually form any adjacent quadgrams, so we can’t use this method.

The best criterion we have is ‘how similar does this distribution look to English?’. This can be very effectively represented by getting a probability distribution from the sub-text (this is what all of the ‘norm_dist’ stuff was about), and then for each letter taking the difference between its probability in standard english, its probability in the text, and squaring it, taking the sum of these squares. The squaring step is useful as it ensures only the magnitude of a difference is considered - ie negative differences are still differences. This could also have been done using the abs function, but squaring has the added advantage of being ‘harsher’ for larger differences, and less harsh for smaller differences. Using this metric, we can determine each letter of the keyword. The following are functions from ‘UAttack.pas’:

```

1 function normalise(d: dist): norm_dist;
2 var
3     total: integer = 0;
4     i: integer;
5     dist_i: char;
6 begin
7     for i in d do
8         total := total + i;
9     for dist_i := 'A' to 'Z' do
10        normalise[dist_i] := d[dist_i] / total;
11 end;
12
13 function fitness(pt_dist: norm_dist): real;
14 var
15     i: char;
16 begin
17     fitness := 0;
18     for i := 'A' to 'Z' do
19         fitness := fitness + power(pt_dist[i] - eng_dist[i], 2);
20 end;
```

Listing 17: Fitness library functions

And here is the calling program:

```

1 {$MODE OBJFPC}
2
3 program PCrack;
4
5 uses UAttack, UEncrypt, UFiles, Sysutils;
6
7 function likely_caes(pt: ansistring): char;
8 var
9     c, best_c: char;
10    f, best_f: real;
11 begin
12     best_f := -999;
13     for c := 'A' to 'Z' do begin
14         f := -fitness(normalise(get_dist(un_vigenere(pt, c))));
15         if f > best_f then begin
16             best_f := f;
17             best_c := c;
18         end;
19     end;
20     likely_caes := best_c;
21 end;
22
23 function likely_vig(pt: ansistring; keyl: integer): string;
24 var
```

```

25     strpd: ansistring;
26     start: integer;
27 begin
28     strpd := clean(pt);
29     likely_vig := '';
30     for start := 0 to keyl - 1 do
31         likely_vig := likely_vig + likely_caes(get_interval(strpd, start, keyl));
32     end;
33
34     var
35         vg_key: string;
36         pt: ansistring;
37         keyl: integer;
38     begin
39         if paramcount > 0 then
40             try
41                 keyl := strtoint(paramstr(1));
42             except
43                 on EConvertError do begin
44                     writeln(stderr, 'Invalid key length ', paramstr(1));
45                     exit;
46                 end;
47             end
48         else begin
49             writeln(stderr, 'Key length required; see bin/IOC');
50             exit;
51         end;
52         pt := read_file(input);
53         vg_key := likely_vig(pt, keyl);
54         writeln(output, 'key is ', vg_key);
55         writeln(output, 'resulting pt: ', un_vigenere(pt, vg_key));
56     end.

```

Listing 18: Keyword-cracking program (PCrack.pas)

This first actually defines a function ‘likely_caes’ to crack a caesar cipher, determining a single letter. This function is then called from ‘likely_vig’ on each caesar-encrypted subtext to find each letter of the keyword (this of course also works on a caesar cipher, where the length of the keyword is 1).

This also reads the length of the keyword from command-line arguments, which the user should have determined using the IOC program, using various failsafes.

Cracking a cipher

Now I can put it to use. I will use challenge 4b from the cipher challenge as a demonstration. Cracking it takes literally two commands now, each running in a couple of ms:

```

1 $ cat ~/programmeren/cipher_tools/src/samples/4b.txt | bin/IOC
2 1 (0.0429) _____
3 2 (0.0429) _____
4 3 (0.0428) _____
5 4 (0.0428) _____
6 5 (0.0427) _____
7 6 (0.0428) _____
8 7 (0.0427) _____
9 8 (0.0427) _____
10 9 (0.0429) _____
11 10 (0.0426) _____
12 11 (0.0433) _____
13 12 (0.0426) _____
14 13 (0.0688) _____
15 14 (0.0427) _____
16 15 (0.0427) _____
17 16 (0.0427) _____
18 17 (0.0428) _____
19 18 (0.0431) _____
20 19 (0.0429) _____
21 20 (0.0421) _____
22 $ cat ~/programmeren/cipher_tools/src/samples/4b.txt | bin/Crack 13
23 key is ARCANAIMPERII
24 resulting pt: OVER THE YEARS THE HEROIC ROLE OF AGRICOLA AT WATLING STREET...

```

Listing 19: Cracking 4b

Source

All involved files, including this L^AT_EX document, can be found at <https://github.com/elterminad0r/encryption>. An interested reader may also wish to explore https://github.com/elterminad0r/cipher_tools, the repository with all the code I produced for the cipher challenge this year. It includes automated crackers for the Hill cipher, autokey cipher and affine Vigenère ciphers by brute force, a framework to allow a user to make substitutions to a text in aid of the cracking of a generic substitution cipher, and crackers for substitution ciphers, playfair ciphers and bifid ciphers written in C using simulated annealing (including quadgram scoring). It also features a slightly derelict Python script that tries to split punctuation-stripped text into words using a prefix tree.