

The application of noisy-channel coding techniques to synthetic DNA barcoding

Name: Izaak van Dongen
EP Mentor: Nicolle Mcnaughton
Tutor: Paul Ingham
Candidate No: 6659

July 15, 2018

Contents

1	Definitions and conventions	2
2	Introduction	5
3	Applications	6
4	The Hamming distance	6
5	Parity codes	7
5.1	Parity	7
5.2	Simple example	7
5.3	Limitations	8
6	The Hamming code	9
6.1	Description	9
6.2	Implementation	10
6.3	Adaption	12
6.3.1	Naïve approach	12
6.3.2	Clever approach	12
7	The Hadamard code	13
7.1	Description	13
7.2	Implementation	14
8	Unit tests	16
9	Miscellaneous listings	18
10	Source	22
	References	22

List of Listings

1	Example of a code listing with comments	5
2	Binary Hamming code in Python	12
3	Hadamard matrix generation	16
4	Unit tests for encode_hamming	17
5	Unit tests for hadamard_matrix	18
6	Generating Hamming coverage indices	18
7	Generating binary table	18
8	Hamming index coverage	20
9	Hadamard visualisation (uses code in 10)	21
10	Template for Hadamard graphic	22

List of Figures

1	Example of a CLI interface	3
2	Illustration of CRISPR editing DNA	6
3	Asymptotic performance of log and $\sqrt{\cdot}$	8
4	Index coverage of Hamming parity bits	10
5	Visualisations of “2 ⁿ ” Hadamard matrices	14

List of Tables

1	Parity coverage in a Hamming code	9
2	Indices covered by each parity bit shown in binary	10

1 Definitions ¹ and conventions

My project will involve a lot of programming. Because of the tendency for programmers to use field-specific vocabulary/jargon, I have provided the following definitions of various words I might use when talking about programs I write.

Definition:

A **program** or **script**, but not to be confused with) is a series of instructions that the computer follows in order to complete a task. They are normally created and shown in the form of text. Roughly speaking, each line of text corresponds to one instruction for the computer to follow. A program in this form is written in a programming language.

Definition:

A program normally also permits **input**, **arguments** or **parameters**. These allow a user of the program to feed it some data or starting instructions. This is extremely useful as it means that the program doesn’t just do the same thing each time, but performs its task on the data or information requested by the user.

¹Please note that these definitions have been somewhat simplified. It would be impractical and outside of the scope of this dissertation to deliver a full briefing of the field of computer science.

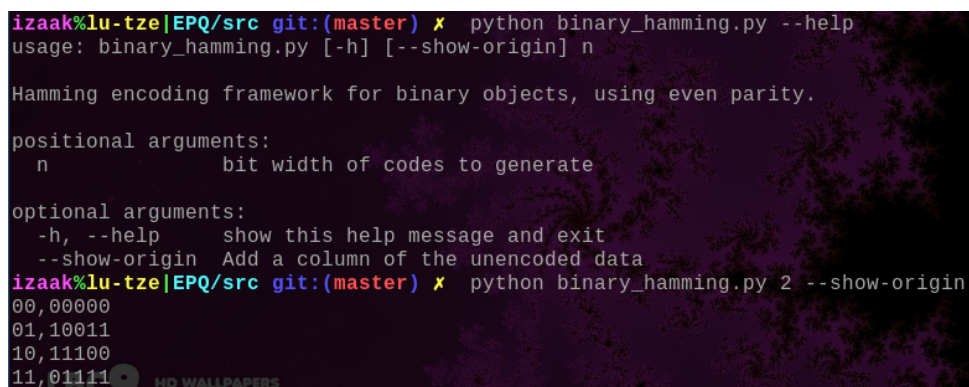
Definition:

The word **code** sits kind of annoyingly here. Code may both refer to programming instructions or systems of symbols representing data. When referring to programs, code is an uncountable noun, eg “I wrote some code”. When referring to symbolic systems, it is countable, so you might say “I generated some codes, using the code I wrote earlier.”

Definition:

The programs I produce will be CLI-based. That is, they will use a **Command Line Interface**. This means that the user interacts with the program solely through text. Most applications nowadays use a graphical interface, but a CLI has many benefits, including but not limited to being easier to develop, being faster due to less overhead, CLI programs can interact with each other more easily as it provides a standard interface, and they are more portable. The downside of course is that it required more expertise to use.

An example of how I use the CLI to interact with a program is given in Figure 1.



```
izaak%lu-tze|EPQ/src git:(master) x python binary_hamming.py --help
usage: binary_hamming.py [-h] [--show-origin] n

Hamming encoding framework for binary objects, using even parity.

positional arguments:
  n                    bit width of codes to generate

optional arguments:
  -h, --help            show this help message and exit
  --show-origin          Add a column of the unencoded data
izaak%lu-tze|EPQ/src git:(master) x python binary_hamming.py 2 --show-origin
00,00000
01,10011
10,11100
11,01111
```

Figure 1: Example of a CLI interface

Definition:

A **programming language** is a defined language that both the computer and the programmer understand. The programming language I’m using for this dissertation is Python. To produce some of by graphics, I’m using Postscript. On a more meta-level this dissertation itself has been produced with \LaTeX , and I have worked on it on a GNU/Linux system, making use of the shell language Zsh, and other tools like the language ‘Make’. None of these are particularly important to understanding the outcome or goals of this dissertation.

Definition:

A **comment** is a section of a program which has been specially marked to be ignored by the computer. These are used by humans to add documentation or clarification to code, and are written in natural language. In the two languages I’m using, comments are marked by a preceding % or #.

Definition:

A **docstring** is also text within a program. It is similar to a comment, but there are some technical and semantic differences. Most importantly, a docstring is intended to describe the function of some unit of the code. This might be one function, one class, or of the whole program. A docstring is enclosed by three consecutive double quotes: `"""`. By convention, docstrings are written in the imperative mood [9].

Definition:

A **function** is a part of the code that acts out one specific task. These are useful as they can be reused, making code more maintainable (only one part must be modified), shorter and easier to test. Functions in Python are generally preceded by `def` `function_name(parameters):`. Here, 'parameters' defines which inputs the function expects. This allows the function to perform the same task on different data. All of the code within the function will then be indented one level.

Definition:

One program can reuse a function from another program. To do this, the requesting program must **import** the function from the defining program. Most language also have a standard library of useful functions, so some imports will be from the language itself, rather than another file within this EP.

Definition:

A **unit test** is a way to test that a program is functioning correctly. It works by running a number of tests against each *unit* of code. These units are normally functions. By making sure that each function works, we can be reasonably confident in how robust our code is. A unit test is itself also a program.

Definition:

A **matrix** is a rectangular array of numbers. They are denoted like: $M = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 5 \end{pmatrix}$. Matrices have **rows** and **columns**, which may have certain properties desirable for barcodes.

For my EP I have written several programs. Almost all of these are included within the dissertation or some other component as code listings.. When appropriate, I have added comments and docstrings to these, to explain (commentate) their function.

In the code listings provided, comments should be highlighted *#like this*, and docstrings `"""like this"""`. All of this is demonstrated in listing 1.

```

1 #!/usr/bin/env python3
2 ^ This part is a shebang. It tells the computer what language I'm using
3
4 """
5 This part is a docstring and explains what all of the code does

```

```
6  """
7
8  import pprint
9  import json
10
11 def this_is(some: "code") -> {"th": at}:
12     """
13     This is another docstring, but in a function, to describe what the function
14     does.
15     """
16     for x in y:
17         pprint(json.loads('{"a": 2}'))
18         # This part is a comment. It explains what the following line of code
19         # does, because it is particularly interesting/complicated.
20         does(stuff)
21
22         this is {a: really(long-line(of + code, [that, goes], off(the_edge),
23             ↪ resulting in "a little arrow")}
24
25 # <- this is a comment pointing out the line numbers
```

Listing 1: Example of a code listing with comments

2 Introduction

The premise of this project is to investigate the different types of error-correcting codes, and how these might be applied to DNA barcoding. DNA barcoding is the assignment of ‘barcodes’ to subsequences of synthesised DNA for the purposes of identification. This means that in future, when you look at that same subsequence, that should let you ‘ID’ the whole DNA sequence. However, as we all know, DNA is subject to mutation, so an ideal barcode should still be identifiable after some number of mutations. This is where error-correcting codes come in.

The challenge in this comes from the fact that most error-correcting codes are designed in base-2 (binary) whereas DNA strings are fundamentally base-4 (quaternary). The applicability of this project is that in oligonucleotide synthesis, some samples may need to be identified later on using a subsection of the sample (a barcode). These could just be linearly assigned codes, but this would leave them very susceptible to mutation.

Here is an example: say that we’re given a barcode of length four, to encode two different samples. If we worked methodically up from the bottom (using the ordering ACGT - orderings will be discussed further later on) we might end up with the codes AAAA and AAAC. However, either string would only require a single mutation (where we say a mutation is the changing of a single base) to become identical to the other one. Therefore, in this case, it would clearly be far more optimal to make a choice like, for example, AAAA and CCCC (although this leaves you with only two different codes).

We have kind of glossed over how we in this case formally represent DNA and mutations, but we will get to that.

There are also a number of parameters to the problem, and as they change the problem becomes increasingly hard:

- What if the barcode size changes?

- What if we want more codes than two?
- What if we anticipate many mutations?

All of these will be further explored in this dissertation.

3 Applications²

You may well have heard of CRISPR/Cas9, depending on which news sources you follow. CRISPR stands for “Clustered Regularly Interspaced Short Palindromic Repeats”, which is a family of DNA sequences in bacteria [4, 13, 26]. Cas9 is a specific CRISPR Associated Protein [25]. Together, they provide a mechanism for the editing of the DNA of an organism [7], by using a specific kind of virus. This mechanism is shown in figure 2.

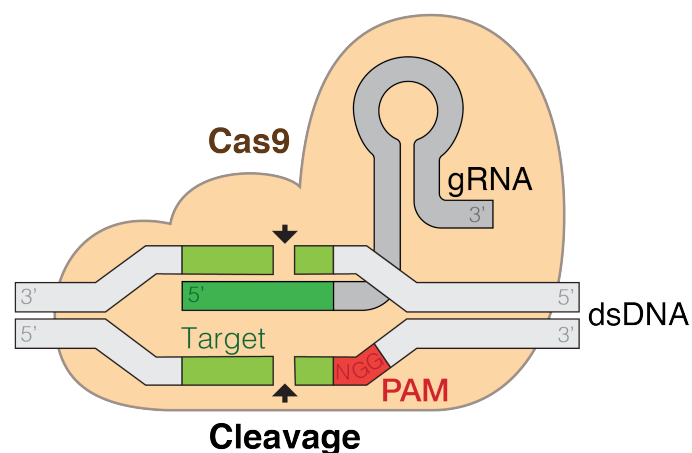


Figure 2: Illustration of CRISPR editing DNA

This mechanism requires “guide RNA (gRNA)”. This RNA may sometimes have to be synthesised on a large scale, and this is where the barcodes come in. If a large “batch” of gRNA contains RNA to be used for different experiments, this could then be barcoded in order to identify the RNA.

This is just one example of what could be done - there are many other applications [27].

4 The Hamming distance

The Hamming distance is a measure of “string distance”. String distance is a way to define how different two strings are. Coding-theoretically, this can be used to quantify the amount that a string has been changed by transmission (or an oligonucleotide has been mutated).

The Hamming distance between any two equally long strings S and R is given by the number of characters at identical position that differ. For example, if we let d_H denote Hamming distance, the distances

$$\begin{aligned} d_H(S, R) &= 1 \\ d_H(S, T) &= 2 \\ \text{where } S &= abcde \end{aligned}$$

²This section gives a brief overview of possible applications - be aware that this is somewhat abbreviated as this topic is not the focus of this dissertation.

$$R = \text{abcfe}$$

$$T = \text{axcze}$$

Note that for any S , $d(S, S) = 0$. This means that there is no “distance” from a string to itself.

In terms of DNA, the Hamming distance can be used to determine the number of bases that have mutated.

5 Parity codes

5.1 Parity

The insertion of “parity bits” is a common practice in basic encoding. Parity refers to the “oddness” or “evenness” of some data [16].

Commonly, this is determined by making sure that the sum of the data is $0 \bmod 2$. This is also known as “even parity”, and means that the total number of 1-bits in the data should be even. This means that the data is known to have been corrupted if its sum is odd.

The reason that I am using such an elaborate definition with sums and modulus is that I know I want to adapt my codes to base 4, and by using such a general definition, this is done quite easily.

5.2 Simple example

A simple but inefficient parity encoding scheme is a column/row wise encoding [10]. Take the slightly contrived data string “01000001010100”. This is very tangentially related to DNA - it’s the 8-bit ASCII [32] representation of the string “AT”, generated by the Python: `"".join(bin(ord(c))[2:].rjust(8, "0") for c in "AT")`³.

The string is then arranged in a square like so:

0	1	0	0
0	0	0	1
0	1	0	1
0	1	0	0

An extra row and column, including an extra corner piece is appended like so:

0	1	0	0	1
0	0	0	1	1
0	1	0	1	0
0	1	0	0	1
0	1	0	0	1

Each of the extra bits documents the parity of its row. Using a scheme like this, a single corrupted bit can be detected, and corrected. For example, the bit at (3, 4) may have flipped like so:

0	1	0	0	1
0	0	0	1	1
0	1	0	1	0
0	1	1	0	1
0	1	0	0	1

³Sometimes I will include some ‘meta-code’ that was used to generate a table or other data. Especially the smaller samples won’t be as extensively commented as I don’t consider these core programs - they are a kind of shortcut from writing the table out by hand.

Someone wishing to correct this error can check the parity of each column, compared with its parity bit. They can do the same for each row. Assuming one error has occurred, the point where the incorrect row and column cross is to be flipped back. In this case, the third column doesn't add up, and the fourth row doesn't add up, leading to the faulty bit. It is worth noting that this also works to correct errors in the parity bits, due to the extra corner bit. If only the extra corner bit seems to be wrong, it is the one that has flipped.

5.3 Limitations

This particular scheme is in a sense quite inefficient. At the most optimal configuration, it uses on the order of $2\sqrt{n}$ parity bits, where n is the number of bits in the message, in order to achieve 1 correction. This can be proven as follows:

Assume n to be highly divisible. Let p denote the number of parity bits, and x denote the length of a row. We then have,

$$\begin{aligned}
 p &= \frac{n}{x} + x \\
 \Rightarrow \frac{dp}{dx} &= 1 - \frac{n}{x^2} = 0 \text{ (as } p \text{ is a minimum)} \\
 \Rightarrow 1 &= \frac{n}{x^2} \\
 \Rightarrow x^2 &= n \\
 \Rightarrow x &= \sqrt{n} \\
 \Rightarrow p &= \frac{n}{\sqrt{n}} + \sqrt{n} = \sqrt{n} + \sqrt{n} \\
 &= 2\sqrt{n}
 \end{aligned}$$

This is quite a poor asymptotic performance - as the number of data bits grows larger, the number of parity bits required grows relatively fast. In the next section, I describe a similar code that uses only $\log_2 n$. In figure 3 is a quick plot comparing the two functions.

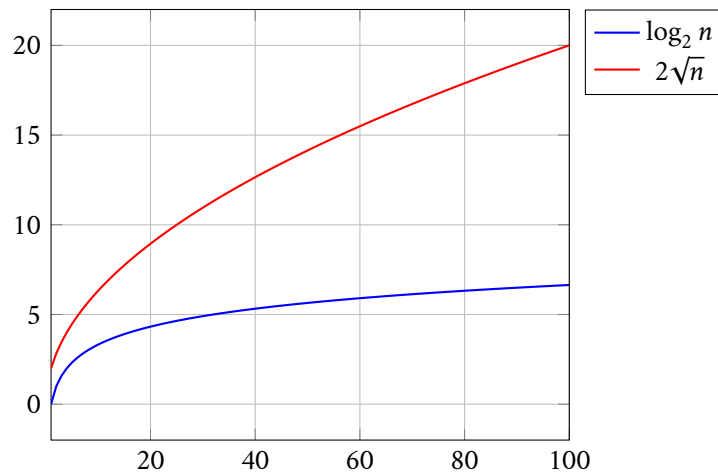


Figure 3: Asymptotic performance of log and $\sqrt{}$

As you can see, as n increases the relative performance of the row-column approach degrades significantly.

6 The Hamming code

6.1 Description

The Hamming code instead places a parity bit at each index that is a power of two, where we number indices starting from 1. Therefore, our previous data string gains parity bits in this configuration: **100110000001010010100** (the 1st, 2nd, 4th, 8th and 16th bits are used for parity).

The way the parity “coverage” works is shown in table 1. I have included indices up to 31. This is because that is the longest encodable string with only five parity bits (afterwards, we have to add a parity bit at 32). Of course, a shorter code word can always also be encoded by just acting as if each index that is out of range stores a 0.

Parity index		Covered indices														
1	3	5	7	9	11	13	15	17	19	21	23	25	27	29	31	
2	3	6	7	10	11	14	15	18	19	22	23	26	27	30	31	
4	5	6	7	12	13	14	15	20	21	22	23	28	29	30	31	
8	9	10	11	12	13	14	15	24	25	26	27	28	29	30	31	
16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	

Table 1: Parity coverage in a Hamming code

These are very deliberately chosen indices. In fact, this table was generated by a short code snippet that can be found in listing 6.

The way that it works is by considering the value of the parity index in binary. For example, $4_{10} = 100_2$. As they are powers of two, they will always be of the form ‘10*’ (a one followed by 0 or more zeroes).

The code in listings 7 and 8 generates two visualisations, which I find helpful. The first is table 2, which is similar to table 1, but transposed so each column corresponds to a parity bit, and each index is written in binary:

Parity index	00001	00010	00100	01000	10000
Coverage	00011	00011	00101	01001	10001
	00101	00110	00110	01010	10010
	00111	00111	00111	01011	10011
	01001	01010	01100	01100	10100
	01011	01011	01101	01101	10101
	01101	01110	01110	01110	10110
	01111	01111	01111	01111	10111
	10001	10010	10100	11000	11000
	10011	10011	10101	11001	11001
	10101	10110	10110	11010	11010
	10111	10111	10111	11011	11011
	11001	11010	11100	11100	11100
	11011	11011	11101	11101	11101
	11101	11110	11110	11110	11110
	11111	11111	11111	11111	11111

Table 2: Indices covered by each parity bit shown in binary

The second is shown in figure 4. It represents each covered bit as a filled in square, and each non-covered bit as an empty square, so the whole codeword is shown in every row.

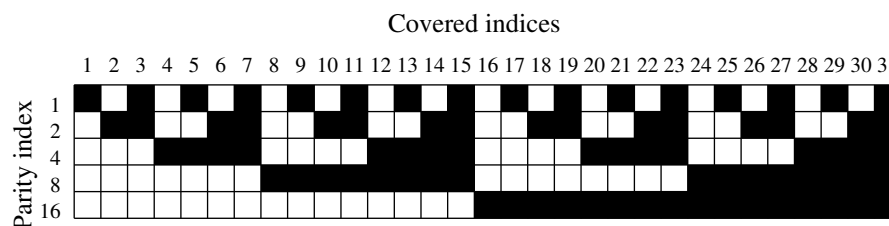


Figure 4: Index coverage of Hamming parity bits

Now, when decoding a Hamming-encoded message, you check each parity bit for errors, as normal. You then add the index of each parity bit that shows an error to find the index of the corrupted bit. This bit can then be flipped. This can only correct one error, but it is a very effective way to do so - in fact, if you only want to correct one error, the Hamming code is optimal [11].

Of course, the Hamming code is not explicitly a way to generate barcodes - it is rather an encoding. However, error-correcting encodings can simply be applied to all possible strings to generate a set of barcodes. As we are using binary for now, you could generate the code corresponding to each of the 16 binary strings from 0000 to 1111. This gives you 16 barcodes of length 7 each, which can be decoded to retrieve the unique identifier from 0 to 15 that was associated with the DNA.

Listing 2 shows the program that implements a binary Hamming encoding. It is tested by the code in listing 4.

6.2 Implementation

```

1  #!/usr/bin/env python3
2
3  """

```

```
4 Hamming encoding framework for binary objects, using even parity.
5 """
6
7 # imports the "argparse" library, which is used to interpret the parameters
8 # given to the program by the user
9 import argparse
10
11 # imports several functions that can be used to manipulate sequences - in this
12 # case sequences of bits
13 from itertools import count, takewhile, product
14
15 def get_args():
16     """
17     Use argparse to interpret parameters:
18     - the number of bits the user wants to encode
19     - if the unencoded data should be shown
20     - what base should be used to calculate the parity?
21     """
22     parser = argparse.ArgumentParser(description=__doc__)
23     parser.add_argument("n", type=int, help="bit width of codes to generate")
24     parser.add_argument("--show-origin", action="store_true",
25                         help="Add a column of the unencoded data")
26     parser.add_argument("--base", type=int, default=2,
27                         help="base to generate codes for")
28     return parser.parse_args()
29
30 def powers_to(n=float("inf")):
31     """
32     Get all of the powers of 2 up to a given upper limit. Uses count() to
33     produce the set of natural numbers (0, 1, 2, 3..) and takewhile() to keep
34     taking powers of 2 until they exceed the limit. By default, produces all
35     powers of 2.
36     """
37     return takewhile(lambda x: x < n, (1 << i for i in count()))
38
39 def matching_indices(power, l):
40     """
41     Generate the particular indices covered by a parity bit.
42     """
43     return (i for pstart in range(power - 1, l, power << 1)
44             for i in range(pstart, min(l, pstart + power)))
45
46 def hamming_encode(bin_stream, base):
47     """
48     Perform Hamming encoding of a series of bits.
49     """
50     power = 1
51     out = []
52
53     # first fill in each index to be used for parity with [False]
54     for bit in bin_stream:
```

```

55     while len(out) + 1 == power:
56         power <<= 1
57         out.append(False)
58     out.append(bit)
59
60     # then go through each parity index and set it to the residue of the sum of
61     # its data bits modulo (base)
62     for power in powers_to(len(out)):
63         out[power - 1] = sum(-out[i] for i in
64                               matching_indices(power, len(out))) % base
65     return out
66
67 # if the program is called directly, generate as many codes as the user wants.
68 if __name__ == "__main__":
69     # uses the previous function to get the user's input
70     args = get_args()
71     # this generates all possible binary strings of length n, by taking the
72     # cartesian product {0..base}^n
73     for code in product(range(args.base), repeat=args.n):
74         # if we also want to show the data bits, print them to the screen
75         if args.show_origin:
76             print("".join(map(str, code)), end=",")
77         # print the encoded data to the screen
78         print("".join(map(str, hamming_encode(code, args.base))))

```

Listing 2: Binary Hamming code in Python

6.3 Adaption

Unfortunately, this all produces binary codes, as this is more of a ‘fundamental’ base in the world of computers.. As is, this is of no use because DNA strings are fundamentally base-4.

6.3.1 Naïve approach

We are quite fortunate in that 4 is a power of 2. This means that we can directly translate a base-4 string to binary, in the case of DNA perhaps by mapping ACGT to 00 01 10 11.

This specific ordering is quite useful as each base pair A-T and C-G is a set of additive inverses mod 4, or “number bonds to 4”.

In any case, by writing a short program to perform this mapping, it is quite trivially done.

6.3.2 Clever approach

It is also possible to directly translate

7 The Hadamard code

7.1 Description

The Hadamard code is based on Hadamard matrices. A Hadamard matrix is a matrix such that each pair of rows represent a pair of orthogonal vectors [12]. Practically, this means that each row has a Hamming distance of at least half of its length from each other row. This is a much stronger encoding than the Hamming code, so may be much more resistant to mutations. However, as a natural side effect of this, Hadamard codes are longer and more sparse.

A basic generation scheme for a sequence of Hadamard matrices H_{2^n} is given by Sylvester's construction [23]:

$$\begin{aligned} H_1 &= (1) \\ H_2 &= \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \\ H_{2^{n+1}} &= \begin{pmatrix} H_{2^n} & H_{2^n} \\ H_{2^n} & -H_{2^n} \end{pmatrix} = H_2 \otimes H_{2^n} \end{aligned}$$

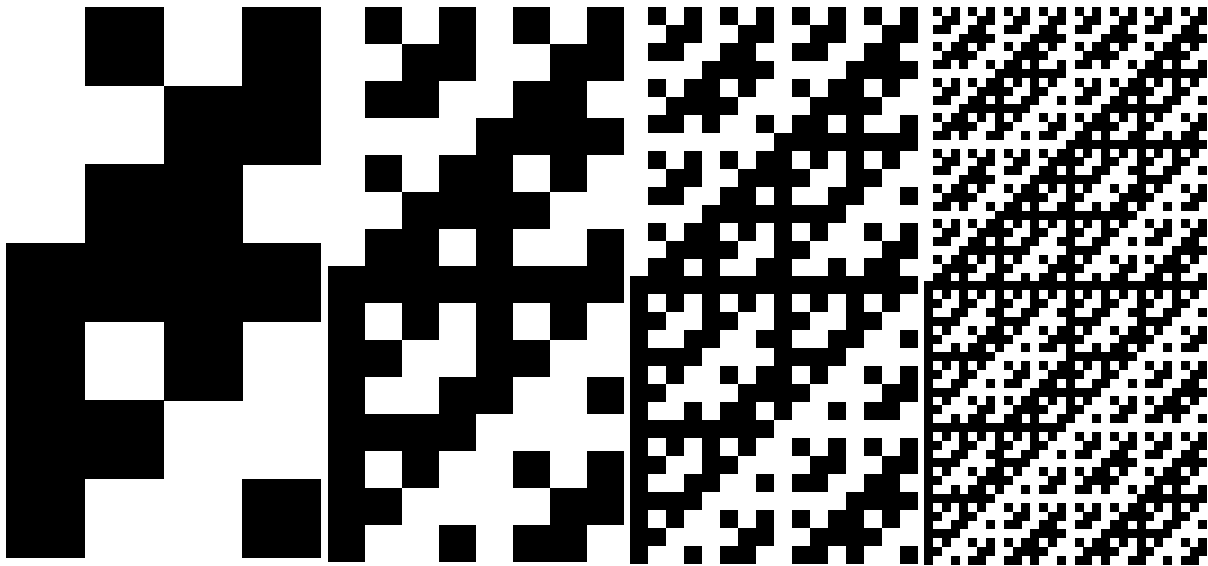
ie the next matrix is

$$H_4 = \begin{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} & \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \\ \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} & -\begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{pmatrix}$$

Note this construction's use of the $(-)$ operator, together with the set $\{1, -1\}$ as it's "alphabet". This is necessary for the satisfaction of the "orthogonality" definition of Hadamard matrices, but, in fact, we need only a slightly weaker property, which is half of the items being different. I have therefore slightly modified this construction to use the Boolean alphabet $\{0, 1\}$, and use the Boolean logical not (\neg) [15]. This is more suited to this application of these codes, as 0 and 1 are the fundamental numbers that computers work with (known as "bits"), so this makes the implementation more natural and elegant, as most programming languages have built in support for Boolean logic.

The code implementing Sylvester's construction can be found in listing 3.

A visualisation of the Hadamard matrix is provided by the Python script in listing 9. It displays the matrix as a grid, where each '1' is filled in, as in figure 5.

Figure 5: Visualisations of “ 2^n ” Hadamard matrices

7.2 Implementation

```

1  #!/usr/bin/env python3
2
3  """
4  Generating a (binary) Hadamard matrix
5  """
6
7  # imports the "system" library, which is used to write error messages to
8  # sys.stderr
9  import sys
10
11 # imports the "time" library, which is used to show diagnostic timing
12 # information (see get_matrix)
13 import time
14
15 # imports the "argparse" library, which is used to interpret the parameters
16 # given to the program by the user
17 import argparse
18
19 def get_args():
20     """
21     Use argparse to get:
22     - The numbers of Hadamard iterations to perform
23     - A file to store the resulting matrix in (they can get quite large)
24     - Whether or not to show diagnostic timing
25     - How to format the matrix
26     """
27     parser = argparse.ArgumentParser(description=__doc__)
28     parser.add_argument("iterations", type=int,
29                          help="Number of iterations to perform on matrix")
30     parser.add_argument("--dump", type=argparse.FileType("w"), default="-",

```

```

31         help="File to write Hadamard matrix to")
32     parser.add_argument("--verbose", action="store_true",
33                         help="Write diagnostic information to stderr")
34     parser.add_argument("--pretty", action="store_true",
35                         help="Use visual block character to display 1")
36     return parser.parse_args()
37
38 def prettify(had_mat, t_char="x", f_char=" "):
39     """
40     Format a hadamard matrix nicely, using 'x' to represent a 1, by default.
41     """
42     return "\n".join("".join(t_char if i else f_char for i in row)
43                      for row in had_mat)
44
45 def hadamard_iterate(mat):
46     """
47     Perform a Hadamard iteration on a matrix.
48     """
49     for r_ind in range(len(mat)):
50         # add the current row the the end of the matrix, duplicated twice
51         mat.append(mat[r_ind] * 2)
52         # add the row's own inverse to its end
53         mat[r_ind].extend([not i for i in mat[r_ind]])
54
55 def get_matrix(iterations, verbose=False):
56     """
57     Generate a full Hadamard matrix given number of iterations
58     """
59     # the start time
60     start = time.time()
61     # the initial matrix
62     had_mat = [[1]]
63     # iterate the appropriate number of times
64     for i in range(iterations):
65         hadamard_iterate(had_mat)
66         # if the user wants to know, provide diagnostic information
67         if verbose:
68             sys.stderr.write("iteration {} successful at {:.3f}s"
69                             .format(i, time.time() - start))
70     # Finally, append -M to the bottom of M
71     for r_ind in range(len(had_mat)):
72         had_mat.append([not i for i in had_mat[r_ind]])
73     return had_mat
74
75 # if the script is called directly, generate a matrix and format and write it
76 # as specified
77 if __name__ == "__main__":
78     args = get_args()
79     display_chars = "10"
80     if args.pretty:
81         display_chars = "\u2588\u2588", " "

```

```

82     print(prettify(get_matrix(args.iterations, args.verbose), *display_chars),
83           file=args.dump)

```

Listing 3: Hadamard matrix generation

8 Unit tests

For all the core programs, I have also written unit tests. These aim to verify that the program works by presenting a number of test cases, and seeing if the program produces the correct output. These both help to ensure correct behaviour, and can serve as a more practical reference of how I expect functions to behave.

```

1  #!/usr/bin/env python3
2
3  """
4  Unit tests for encode_hamming.py.
5
6  This is a testing program, that runs a series of test cases to verify that my
7  code works.
8  """
9
10 import unittest
11
12 from encode_hamming import powers_to, hamming_encode, matching_indices
13
14 class BinaryHammingTestCase(unittest.TestCase):
15     # testing the "powers_to" function
16     def test_powers_to(self):
17         self.assertEqual(list(powers_to(0)), [])
18         self.assertEqual(list(powers_to(1)), [])
19         self.assertEqual(list(powers_to(2)), [1])
20         self.assertEqual(list(powers_to(4)), [1, 2])
21         self.assertEqual(list(powers_to(5)), [1, 2, 4])
22         self.assertEqual(list(powers_to(13)), [1, 2, 4, 8])
23
24     # testing the "hamming_encode" function
25     def test_hamming_encode(self):
26         self.assertEqual(hamming_encode([1, 0, 1, 1], 2), [0, 1, 1, 0, 0, 1, 1])
27         self.assertEqual(hamming_encode(
28             [0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0], 2),
29             [1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0])
30
31     # testing the "matching_indices" function
32     def matching_indices(self):
33         self.assertEqual(matching_indices(1, 7), [1, 3, 5, 7])
34         self.assertEqual(matching_indices(2, 7), [2, 3, 6, 7])
35         self.assertEqual(matching_indices(4, 7), [4, 5, 6, 7])
36         self.assertEqual(matching_indices(4, 1), [])
37         self.assertEqual(matching_indices(1, 1), [1])
38         self.assertEqual(matching_indices(1, 2), [1])

```



```

39         self.assertEqual(matching_indices(4, 5), [4, 5])
40
41 if __name__ == "__main__":
42     unittest.main()

```

Listing 4: Unit tests for encode_hamming

```

1  #!/usr/bin/env python3
2
3  """
4  Unit tests for hadamard_matrix.py.
5
6  This is a testing program, that runs a series of test cases to verify that my
7  code works.
8  """
9
10 import unittest
11
12 from hadamard_matrix import hadamard_iterate, get_matrix
13
14 def make_ints(M):
15     return [[int(i) for i in row] for row in M]
16
17 class HadamardMatrixTestCase(unittest.TestCase):
18     # tests hadamard_iterate
19     def test_hadamard_iterate(self):
20         mat = [[1]]; hadamard_iterate(mat)
21         self.assertEqual(mat,
22             [[1, 0],
23              [1, 1]])
24         mat = [[0]]; hadamard_iterate(mat)
25         self.assertEqual(mat,
26             [[0, 1],
27              [0, 0]])
28         mat = [[1, 0],
29                 [0, 1]]; hadamard_iterate(mat)
30         self.assertEqual(mat,
31             [[1, 0, 0, 1],
32              [0, 1, 1, 0],
33              [1, 0, 1, 0],
34              [0, 1, 0, 1]])
35
36     # tests get_matrix
37     def test_get_matrix(self):
38         self.assertEqual(make_ints(get_matrix(0)),
39             [[1],
40              [0]])
41         self.assertEqual(make_ints(get_matrix(1)),
42             [[1, 0],
43              [1, 1],
44              [0, 1],

```

```

45         [0, 0]])
46
47 if __name__ == "__main__":
48     unittest.main()

```

Listing 5: Unit tests for hadamard_matrix

9 Miscellaneous listings

Below are all the listings that I don't consider to be important to the DNA barcoding part of my dissertation, but have still included as it is material that I have produced for my EP, and illustrate some of the work that has gone into the production of the actual dissertation.

```

1 for p_ind in (1 << pwr for pwr in range(5)):
2     print(r"    {} \\".format(" & ".join(str(i) for i in range(1, 33) if i & p_ind
    ↪ ))))

```

Listing 6: Generating Hamming coverage indices

```

1 def add_color(s, ind):
2     return r"{}\textcolor{blue}{{{}}}".format(s[:ind], s[ind], s[ind+1:])
3
4 table = zip(*[i for i in range(1, 33) if i & p_ind] for p_ind in (1 << pwr for pwr
    ↪ in range(5)))
5 print("\n".join(r"    {} \\".format(" & ".join(r"\texttt{{{}}}".format(add_color(
    ↪ bin(i)[2:].rjust(5, "0"), 4 - sig_ind))
6         for sig_ind, i in enumerate(row)))
7         for row in table))

```

Listing 7: Generating binary table

```

1 %!PS-Adobe-3.0
2
3 /roman {
4     /Times-Roman findfont
5     exch scalefont
6     setfont
7 } def
8
9 /center {
10    /txt exch def
11    /y exch def
12    /x exch def
13
14    txt dup stringwidth pop
15    2 div
16    x exch sub
17    y moveto
18 } def
19
20 /right {

```

```
21     /txt exch def
22     /y exch def
23     /x exch def
24
25     txt dup stringwidth pop
26     x exch sub
27     y moveto
28 } def
29
30 /square {
31     /y exch def
32     /x exch def
33     newpath
34     x y moveto
35     x y 1 add lineto
36     x 1 add y 1 add lineto
37     x 1 add y lineto
38     closepath fill
39 } def
40
41 0.8 roman
42
43 10 dup scale
44 2 0 translate
45 0.05 setlinewidth
46
47 1 1 31 {
48     /ind exch def
49
50     newpath
51     ind 0.5 add 6.5
52     ind 2 string cvs center show
53
54     newpath
55     ind 6 moveto
56     ind 1 lineto
57     stroke
58
59     0 1 4 {
60         /pos exch def
61         /par 2 pos exp cvi def
62         par ind and 0 eq not {
63             ind 5 pos sub square
64         } if
65     } for
66 } for
67
68 0 1 4 {
69     /pos exch def
70     /par 2 pos exp cvi def
71
```

```
72     newpath
73     0.5 5 pos sub
74     par 2 string cvs right show
75
76     newpath
77     1 pos 1 add moveto
78     32 pos 1 add lineto
79     stroke
80 } for
81
82 newpath
83 1 6 moveto
84 32 6 lineto
85 stroke
86
87 1 roman
88
89 newpath
90 16 8 (Covered indices) center show
91
92 gsave
93 newpath
94 -0.6 3 translate
95 90 rotate
96 0 0 (Parity index) center show
97 grestore
98
99 showpage
```

Listing 8: Hamming index coverage

```
1  #!/usr/bin/env python3
2
3  """
4  Generates Postscript file that draws a Hadamard Matrix with filled in boxes.
5  This is a Python script that uses the other Python program to generate a
6  Hadamard matrix, and then inserts it into the premade Postscript template
7  "hadamard_template.ps", which knows how to draw it.
8  """
9
10 # library used to parse the user's arguments. Basically, helps provide an
11 # interface to the program for the user
12 import argparse
13
14 # Regular Expression library. Is used to process text, in "is_ps_comment"
15 import re
16
17 # Operating System Path library. Used to find the location of the "template"
18 # file.
19 import os.path
20
```

```

21 # re-use the code in the "get_matrix" function
22 from hadamard_matrix import get_matrix
23
24 def is_ps_comment(line):
25     """
26     Determine if a line of code is a comment. r"\s*(?:[^\%]|$)" is a "regular
27     expression" that tells the computer to ignore any line that starts with a
28     single % not followed by ! (a comment)
29     """
30     return re.match(r"\s*(?:[^\%]|$)", line)
31
32 # generates full path of template location
33 TEMPLATE_LOCATION = os.path.join(os.path.dirname(__file__),
34                                   "hadamard_template.ps")
35
36 # Load the Postscript template to add data to
37 with open(TEMPLATE_LOCATION, "r") as psfile:
38     PS_SOURCE = "".join(line for line in psfile if not is_ps_comment(line))
39
40 def get_args():
41     """
42     Interpret the user's arguments:
43     - How many iterations should be performed
44     - Where to write the generated Postscript to
45     """
46     parser = argparse.ArgumentParser(description=__doc__)
47     parser.add_argument("iterations", type=int,
48                         help="number of Hadamard iterations")
49     parser.add_argument("--dump", type=argparse.FileType("w"), default="-",
50                         help="file to write generated postscript to")
51     return parser.parse_args()
52
53 # when the program is run
54 if __name__ == "__main__":
55     # get the user's arguments
56     args = get_args()
57     # generate a matrix
58     mat = get_matrix(args.iterations)
59     # insert the matrix in the template and write it to the output file
60     args.dump.write(PS_SOURCE.replace("$HAD_MATRIX",
61     "\n".join("{}").format(" ".join(str(int(i)) for i in row)) for row in mat
61     ↪ ))))

```

Listing 9: Hadamard visualisation (uses code in 10)

```

1 !PS-Adobe-3.0
2
3 % This file is a template used by generate_ham_vis.py
4
5 10 dup scale
6

```

```

7  [
8  % The hadamard matrix is inserted here
9  $HAD_MATRIX
10 ]
11 {
12     % moves "up" by 1 unit for each row
13     0 1 translate
14     gsave
15     {
16         % moves "across" by 1 unit for each square
17         1 0 translate
18         % if the value of the cell in the matrix is 1
19         1 eq {
20             % draw a black square, by making a "path" between the points
21             % (0, 0), (1, 0), (1, 1), (0, 1) and then filling it
22             newpath
23             0 0 moveto
24             1 0 lineto
25             1 1 lineto
26             0 1 lineto
27             closepath fill
28         } if
29         % Do this for all squares in the row
30     } forall
31     grestore
32 % Do this for all rows in the matrix
33 } forall
34
35 % Display this picture
36 showpage

```

Listing 10: Template for Hadamard graphic

10 Source

This document consists of about 3014 words.

All code and source \TeX / \LaTeX files can be found at <https://github.com/elterminad0r/EPQ>.

Unless stated otherwise, all work has been produced by me, including graphics, and source code.

References

- [1] Assmus, E. F. and Key, J. D. [1992], ‘Hadamard matrices and their designs: A coding-theoretic approach’, *Transactions of the American Mathematical Society* **330**(1), 269–293.
URL: <http://www.jstor.org/stable/2154164>
- [2] Baylis, J. [2010], ‘Codes, not ciphers’, *The Mathematical Gazette* **94**(531), 412–425.
URL: <http://www.jstor.org/stable/25759725>

- [3] Bystrykh, L. V. [2012], ‘Generalized DNA barcode design based on Hamming codes’, *PLOS ONE*.
URL: <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0036852>
- [4] Daily Mail [2015], ‘WHAT IS CRISPR-CAS9?’.
URL: <http://www.dailymail.co.uk/sciencetech/fb-5151843/WHAT-CRISPR-CAS9.html>
- [5] del Río, Á. and Rifà, J. [2012], ‘Families of Hadamard $Z_2Z_4Q_8$ -codes’, *CoRR* **abs/1211.5251**.
URL: <http://arxiv.org/abs/1211.5251>
- [6] Ehrenborg, R. [2006], ‘Decoding the Hamming code’, *Math Horizons* **13**(4), 16–17.
URL: <http://www.jstor.org/stable/25678619>
- [7] Finneran, K. [2016], ‘Responding to CRISPR/cas9 [sic]’, *Issues in Science and Technology* **32**(3), 33–34.
URL: <http://www.jstor.org/stable/24727056>
- [8] Golomb, S. W. and Baumert, L. D. [1963], ‘The search for Hadamard matrices’, *The American Mathematical Monthly* **70**(1), 12–17.
URL: <http://www.jstor.org/stable/2312777>
- [9] Goodger, D., van Rossum, G. et al. [2014], Docstring conventions, PEP, Python Software Foundation.
- [10] Guruswami, V. [2010], ‘Introduction to coding theory’.
URL: <http://www.cs.cmu.edu/~venkatg/teaching/codingtheory/notes/notes1.pdf>
- [11] Hamming, R. W. [1950], ‘Error detecting and error correcting codes’, *The Bell System Technical Journal* **26**(2), 147–160.
URL: <http://sb.fluomedia.org/hamming/>
- [12] Hedayat, A. and Wallis, W. D. [1978], ‘Hadamard matrices and their applications’, *The Annals of Statistics* **6**(6), 1184–1238.
URL: <http://www.jstor.org/stable/2958712>
- [13] Horvath, P. and Barrangou, R. [2010], ‘CRISPR/cas, the immune system of bacteria and archaea’, *Science* **327**(5962), 167–170.
URL: <http://www.jstor.org/stable/40508808>
- [14] Hu, W., Kaminski, R., Yang, F., Zhang, Y., Cosentino, L., Li, F., Luo, B., Alvarez-Carbonell, D., Garcia-Mesa, Y., Karn, J., Mo, X. and Khalili, K. [2014], ‘RNA-directed gene editing specifically eradicates latent and prevents new HIV-1 infection’, *Proceedings of the National Academy of Sciences of the United States of America* **111**(31), 11461–11466.
URL: <http://www.jstor.org/stable/23805815>
- [15] Kneale, W. [1956], ‘Boole and the algebra of logic’, *Notes and Records of the Royal Society of London* **12**(1), 53–63.
URL: <http://www.jstor.org/stable/530792>
- [16] Knuth, D., Chapman, R. and Martin, R. [2008], ‘Perfect parity patterns: 11243’, *The American Mathematical Monthly* **115**(7), 668–670.
URL: <http://www.jstor.org/stable/27642574>
- [17] Oztas, E. S. and Siap, I. [2013], ‘Lifted polynomials over F_{16} and their applications to dna codes’, *Filomat* **27**(3), 459–466.
URL: <http://www.jstor.org/stable/24896375>
- [18] Petoukhov, S. V. [2008], The degeneracy of the genetic code and Hadamard matrices, in ‘XX International Congress of Genetics’.
URL: <https://arxiv.org/pdf/0802.3366.pdf>

- [19] Pless, V. [1978], 'Error correcting codes: Practical origins and mathematical implications', *The American Mathematical Monthly* **85**(2), 90–94.
URL: <http://www.jstor.org/stable/2321784>
- [20] Sgarro, A. [1977], 'A fuzzy hamming distance', *Bulletin mathématique de la Société des Sciences Mathématiques de la République Socialiste de Roumanie* **21** (69)(1/2), 137–144.
URL: <http://www.jstor.org/stable/43680301>
- [21] Shannon, C. E. [1948], 'A mathematical theory of communication', *The Bell System Technical Journal* **27**, 379–423, 623–656.
URL: <http://affect-reason-utility.com/1301/4/shannon1948.pdf>
- [22] Spence, E. [1972], 'Hadamard designs', *Proceedings of the American Mathematical Society* **32**(1), 29–31.
URL: <http://www.jstor.org/stable/2038298>
- [23] Sylvester, J. [1867], 'Thoughts on inverse orthogonal matrices, simultaneous signsuccessions, and tessellated pavements in two or more colours, with applications to Newton's rule, ornamental tile-work, and the theory of numbers', *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **34**(232), 461–475.
URL: <https://doi.org/10.1080/14786446708639914>
- [24] Trinh, Q. and Fan, P. [2008], 'Construction of multilevel Hadamard matrices with small alphabet', **44**, 1250–1252.
- [25] Wikipedia [2018a], 'Cas9 — Wikipedia, the free encyclopedia'. [Online; accessed 11 February 2018].
URL: <https://en.wikipedia.org/wiki/Cas9>
- [26] Wikipedia [2018b], 'CRISPR — Wikipedia, the free encyclopedia'. [Online; accessed 10 February 2018].
URL: <https://en.wikipedia.org/wiki/CRISPR>
- [27] Wikipedia [2018c], 'DNA barcoding — Wikipedia, the free encyclopedia'. [Online; accessed 15 June 2018].
URL: https://en.wikipedia.org/wiki/DNA_barcoding
- [28] Wikipedia [2018d], 'Hadamard code — Wikipedia, the free encyclopedia'. [Online; accessed 10 April 2018].
URL: https://en.wikipedia.org/wiki/Hadamard_Code
- [29] Wikipedia [2018e], 'Hadamard matrix — Wikipedia, the free encyclopedia'. [Online; accessed 10 April 2018].
URL: https://en.wikipedia.org/wiki/Hadamard_Matrix
- [30] Wikipedia [2018f], 'Hamming code — Wikipedia, the free encyclopedia'. [Online; accessed 3 April 2018].
URL: https://en.wikipedia.org/wiki/Hamming_Code
- [31] Wikipedia [2018g], 'Hamming(7, 4) — Wikipedia, the free encyclopedia'. [Online; accessed 3 April 2018].
URL: [https://en.wikipedia.org/wiki/Hamming\(7,4\)](https://en.wikipedia.org/wiki/Hamming(7,4))
- [32] X3.4-1963 [1963], American standard code for information interchange, Standard, American Standards Association.
- [33] Yu, Q., Maddah-Ali, M. A. and Avestimehr, A. S. [2017], 'Polynomial codes: an optimal design for high-dimensional coded matrix multiplication', *CoRR* **abs/1705.10464**.
URL: <http://arxiv.org/abs/1705.10464>