

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

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

May 5, 2018

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>The Hamming distance</b>	<b>2</b>
<b>3</b>	<b>Parity codes</b>	<b>2</b>
<b>4</b>	<b>The Hamming code</b>	<b>4</b>
<b>5</b>	<b>Adapting the Hamming code</b>	<b>6</b>
<b>6</b>	<b>The Hadamard code</b>	<b>6</b>
<b>7</b>	<b>Miscellaneous listings</b>	<b>6</b>
<b>8</b>	<b>Source</b>	<b>8</b>
	<b>References</b>	<b>8</b>

## Listings

1	Binary Hamming code in Python . . . . .	5
2	binary_hamming unit tests . . . . .	6
3	Generating Hamming coverage indices . . . . .	6
4	Generating binary table . . . . .	7
5	Graphical index coverage program . . . . .	7

## 1 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. 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.

There have been a few assumptions and glossed over definitions here:

- What constitutes a mutation?
- What is the best way to represent DNA mathematically?

There are also a number of parameters to the problem, and as they change the problem becomes very much nontrivial:

- What if the barcode size changes?
- What if we want more codes than two?
- What if rather than number of codes and barcode size, the parameters are set to barcode size and maximum number of mutations that can occur?

All of these will be further explored in this dissertation.

## 2 The Hamming distance

The Hamming distance is a measure of “string distance”. String distance is a way to define how different two string 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, the distances

$$d(S, R) = 1$$

$$d(S, T) = 2$$

where  $S = \text{abcde}$

$R = \text{abcfe}$

$T = \text{axcze}$

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.

## 3 Parity codes

The insertion of “parity bits” is a common practice in basic encoding. Parity refers to the “oddness” or “evenness” of some data. Commonly, this is determined by the sum of the data modulo 2. For example,

“00101” results in a parity bit of 0, because the sum of all the bits is 2, which has a remainder of 0 when divided by 2 (is equal to 0 mod 2).

A simple but inefficient parity encoding scheme is a column/row wise encoding. Take the slightly contrived data string “0100000101010100”. This is very tangentially related to DNA - it’s the 8-bit ASCII representation of the string “AT”, generated by the Python: `"".join(bin(ord(c))[2:].rjust(8, "0") for c in "AT")`

The string is then split into 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

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 the the extra corner bit. If only the extra corner bit seems to be wrong, it is the one that has flipped.

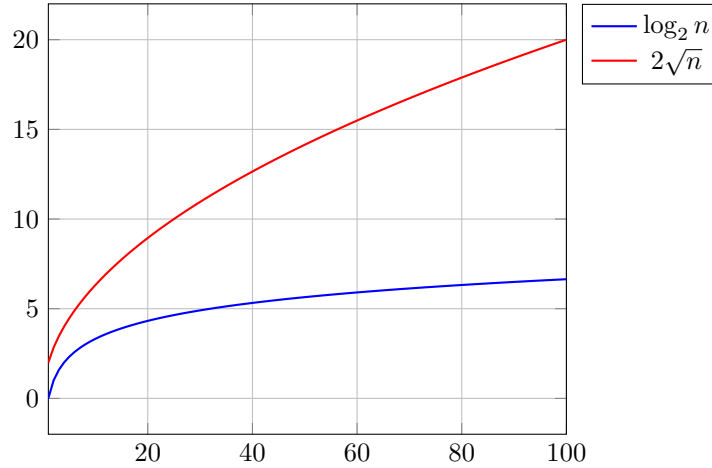
However, this particularly scheme is in a sense quite inefficient. At the most optimal configuration, it uses on the order of  $2\sqrt{n}$  parity bits, where  $n$  is the number of bits in the message, in order to achieve 1 correction. This can be proven as follows:

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,

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

$$= 2\sqrt{n}$$

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  $\log_2 n$ . Here is a quick plot comparing the two functions:



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

## 4 The Hamming code

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 as follows:

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

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 is a 0.

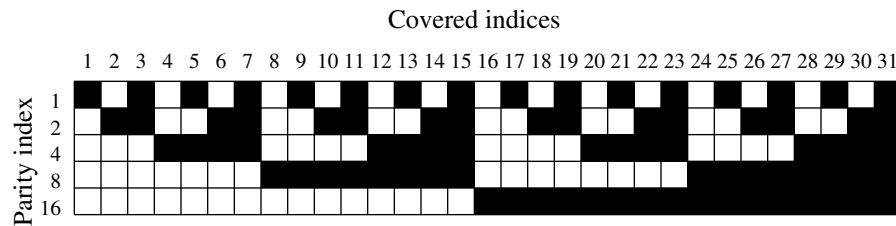
These are very deliberately chosen indices. In fact, this table was generated by a short code snippet that can be found in listing 3.

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 4 and 5 generates the following two visualisations, which I find helpful. The first is a table similar to the first, 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

The second 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:



The script implementing a simple binary Hamming code is as follows:

```

1 #!/usr/bin/env python3
2
3
4 """
5 Hamming encoding framework for binary objects, using even parity.
6 """
7
8 from itertools import count, takewhile
9
10 def powers_to(n):
11     return takewhile(lambda x: x < n, (1 << i for i in count()))
12
13 def matching_indices(power, l):
14     return (i for pstart in range(power - 1, l, power << 1)
15           for i in range(pstart, min(l, pstart + power)))
16
17 def hamming_encode(bin_stream):
18     pwr = 1
19     out = []
20
21     for bit in bin_stream:
22         while len(out) + 1 == pwr:
23             pwr <<= 1
24             out.append(False)

```

```

25         out.append(bit)
26
27     for power in powers_to(len(out)):
28         out[power - 1] = 1 & sum(out[i] for i in matching_indices(power, len(out)))
29     return out

```

Listing 1: Binary Hamming code in Python

This code is accompanied by the following testing scheme:

```

1  """
2  Unit tests for binary_hamming.py
3  """
4
5  import unittest
6
7  from binary_hamming import powers_to, hamming_encode
8
9  class BinaryHammingTestCase(unittest.TestCase):
10     def test_powers_to(self):
11         self.assertEqual(list(powers_to(0)), [])
12         self.assertEqual(list(powers_to(1)), [])
13         self.assertEqual(list(powers_to(2)), [1])
14         self.assertEqual(list(powers_to(4)), [1, 2])
15         self.assertEqual(list(powers_to(5)), [1, 2, 4])
16         self.assertEqual(list(powers_to(13)), [1, 2, 4, 8])
17
18     def test_hamming_encode(self):
19         self.assertEqual(hamming_encode([1, 0, 1, 1]), [0, 1, 1, 0, 0, 1, 1])
20         self.assertEqual(hamming_encode(
21             [0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0]),
22             [1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0])
23
24 if __name__ == "__main__":
25     unittest.main()

```

Listing 2: binary\_hamming unit tests

## 5 Adapting the Hamming code

Unfortunately, this all only operates on binary data, as this is more of a ‘fundamental’ base.. As is, this is of no use because DNA strings are fundamentally base-4.

## 6 The Hadamard code

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. 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

## 7 Miscellaneous listings

```

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 3: 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 4: 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 5: Graphical index coverage program

## 8 Source

This document consists of about 1188 words.

All code and source T<sub>E</sub>X/L<sup>A</sup>T<sub>E</sub>X files can be found at <https://github.com/elterminad0r/EPQ>.

## References

- 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>



- Baylis, J. [2010], ‘Codes, not ciphers’, *The Mathematical Gazette* **94**(531), 412–425.  
**URL:** <http://www.jstor.org/stable/25759725>
- 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>
- 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>
- Ehrenborg, R. [2006], ‘Decoding the hamming code’, *Math Horizons* **13**(4), 16–17.  
**URL:** <http://www.jstor.org/stable/25678619>
- 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>
- Guruswami, V. [2010], ‘Introduction to coding theory’.  
**URL:** <http://www.cs.cmu.edu/~venkatg/teaching/codingtheory/notes/notes1.pdf>
- 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/>
- 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>
- 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>
- 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>
- Petoukhov, S. V. [2008], ‘The degeneracy of the genetic code and hadamard matrices’.  
**URL:** <https://arxiv.org/pdf/0802.3366.pdf>
- 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>
- 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>
- Spence, E. [1972], ‘Hadamard designs’, *Proceedings of the American Mathematical Society* **32**(1), 29–31.  
**URL:** <http://www.jstor.org/stable/2038298>
- Trinh, Q. and Fan, P. [2008], ‘Construction of multilevel hadamard matrices with small alphabet’, **44**, 1250–1252.
- 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>