

The application of noisy-channel coding techniques to DNA barcoding (Early structure/ideas)

Izaak van Dongen

February 1, 2018

Contents

1	Introduction	1
2	The Hamming distance	1
3	Implementing the Hamming code	1
4	Source	2

Listings

1	Binary Hamming code in Python	1
2	binary_hamming unit tests	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. 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.

The Hamming distance

Implementing the Hamming code

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

```
1 #!/usr/bin/env python3
2
3 """
4 Hamming encoding framework for binary objects, using even parity.
5 """
6
7 from itertools import count, takewhile
8
9 def powers_to(n):
10     return takewhile(lambda x: x < n, (1 << i for i in count()))
11
```

```

12 def hamming_encode(bin_stream):
13     pwr = 1
14     out = []
15
16     for bit in bin_stream:
17         while len(out) + 1 == pwr:
18             pwr <<= 1
19             out.append(0)
20             out.append(bit)
21
22     for i in powers_to(len(out)):
23         out[i - 1] = 1 & sum(out[pbit] for pstart in range(i - 1, len(out), i << 1)
24                               for pbit in range(pstart, pstart + i))
25
26     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
21 if __name__ == "__main__":
22     unittest.main()

```

Listing 2: binary_hamming unit tests

Source

References

- [1] Venkatesan Guruswami. Introduction to coding theory. <http://www.cs.cmu.edu/~venkatg/teaching/codingtheory/notes/notes1.pdf>. Accessed: 26/01/2018.
- [2] R. W. Hamming. Error detecting and error correcting codes. *The Bell System Technical Journal*, 26(2):147–160, 1950. Accessed 26/01/2018 via <http://sb.fluomedia.org/hamming/>.
- [3] Leonid V. Bystrykh. Generalized dna barcode design based on hamming codes. *PLOS ONE*, 2012. Accessed 2/02/2018 via <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0036852>.
- [4] C. E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27:379–423, 623–656, 1948. Accessed 2/02/2018 via <http://affect-reason-utility.com/1301/4/shannon1948.pdf>.