



THE UNIVERSITY
of EDINBURGH

Bioinformatics 1

Assignment 1

Henrique Manuel Martins Ferrolho
s1683857 - henrique.ferrolho@gmail.com

October 26, 2016

Contents

1	Question	3
1.1	What is the name of the disease you have selected?	3
1.2	Explain why it is thought there is a genetic basis for this disease.	3
1.3	What is the human name for the gene that is thought to be involved?	3
1.4	Is this gene known by any other names? Whether yes or no, explain how you investigated this.	3
1.5	Is this gene present in a model organism such as the mouse or the fruit fly?	3
2	Question	4
2.1	Now investigate the structure of the gene. Find the gene in a database. Which chromosome is it located on, and at which position along the chromosome?	4
2.2	What is the structure of the gene, how many exons and introns does it have?	4
2.3	Where did you get this sequence from and what was the unique identifier used so that someone else could be sure they were looking at the same sequence?	4
2.4	How long is the transcript, and what proportion is coding? . . .	4
3	Question	5
3.1	Translate your cDNA sequence into protein/amino acid sequence. How many amino acids does your protein contain?	5
3.2	Of the 64 possible codons available, how many are used?	5
3.3	What is the most common amino acid in the protein?	5
3.4	How many codons for this amino acid exist and how often is each used?	5
4	Question	7
5	Question	9
5.1	Sequences 1 and 2 differ slightly. How does the resulting protein differ? Could this have functional implications?	9
5.2	Now use the Needleman Wunsch algorithm to compare sequence 1 to sequences 3 and 4. Use the scoring: match +2, mismatch -1, indel -1. Perform at least one of these on paper (or both if you wish). On paper, use the first three codons only.	9
5.3	Comparing the bare sequences, what can you conclude about the relatedness of the species?	9
5.4	Extra mark	10
6	Sources	11

1 Question

1.1 What is the name of the disease you have selected?

I have selected *Huntington Disease*, also known as *Huntington Chorea*.

HD is a genetic brain disorder which causes jerky movements, emotional problems, and loss of cognition.

1.2 Explain why it is thought there is a genetic basis for this disease.

A DNA segment known as a [CAG trinucleotide repeat](#) has been consistently detected in people who have been subjected to *Molecular Genetic Testing* and suffer from *HD*.

It is believed that an increase of the *CAG segment* length causes *huntingtin* proteins to be longer. Furthermore, these abnormal proteins get split into toxic segments which accumulate in neurons. This compromises neurons' normal behaviour, and might lead to their death.

The manifestation of these events damage areas of the brain, thus originating the symptoms of a person with *HD*.

1.3 What is the human name for the gene that is thought to be involved?

The official symbol is *HTT*. It's official full name is *huntingtin*. Furthermore, it's gene ID in the [NCBI Gene Database](#) is *3064*.

1.4 Is this gene known by any other names? Whether yes or no, explain how you investigated this.

As one can see in the [NCBI huntingtin page](#), the gene is also known as *HD*, and *IT15*.

1.5 Is this gene present in a model organism such as the mouse or the fruit fly?

Yes, it is. The *NCBI* database contains [orthologs of numerous species](#), including [Mus musculus](#) and [Drosophila melanogaster](#).

2 Question

- 2.1** Now investigate the structure of the gene. Find the gene in a database. Which chromosome is it located on, and at which position along the chromosome?

NCBI page: [HTT huntingtin \[Homo sapiens \(human\) \]](#)

It is located in **Chromosome 4 - NC_000004.12**, from position *3074510* to *3243960*.

- 2.2** What is the structure of the gene, how many exons and introns does it have?

It has **67 exons**, and therefore **66 introns** ($67 - 1$).

Download the gene transcript and the coding sequences for your gene - if multiple are listed, choose the main one, at the top of the list.

- 2.3** Where did you get this sequence from and what was the unique identifier used so that someone else could be sure they were looking at the same sequence?

I downloaded both the *Complete Record FASTA* and *Coding Sequences FASTA nucleotide* files from the [NCBI Gene DB](#).

The sequence unique identifier is *NC_000004.12*, from *GRCh38.p7*.

- 2.4** How long is the transcript, and what proportion is coding?

The *Complete Record* file has a length of **169451** nucleotides. The *Coding Sequences* file is **9429** nucleotides long. Therefore, we can say that only **5.56%** of the entire gene correspond to coding sequences:

$$\frac{9429}{169451} = 0.05564$$

3 Question

3.1 Translate your cDNA sequence into protein/amino acid sequence. How many amino acids does your protein contain?

The translation of the *cDNA* contains **3143** amino acids, which makes sense given the *Coding Sequence* is **9429** nucleotides long:

$$9429/3 = 3143$$

3.2 Of the 64 possible codons available, how many are used?

Out of the 64 possible codons available, **62** codons are used. **UAA** and **UAG** are not used.

3.3 What is the most common amino acid in the protein?

The most common amino acid in the *huntingtin* protein is **Leucine** (*Leu/L*) - see figure [1](#).

3.4 How many codons for this amino acid exist and how often is each used?

Leucine has 6 coding amino acids. Below is a table with their respective frequencies in the *huntingtin* protein.

Codon	Frequency
UUA	70
UUG	156
CUU	168
CUC	178
CUA	68
CUG	295

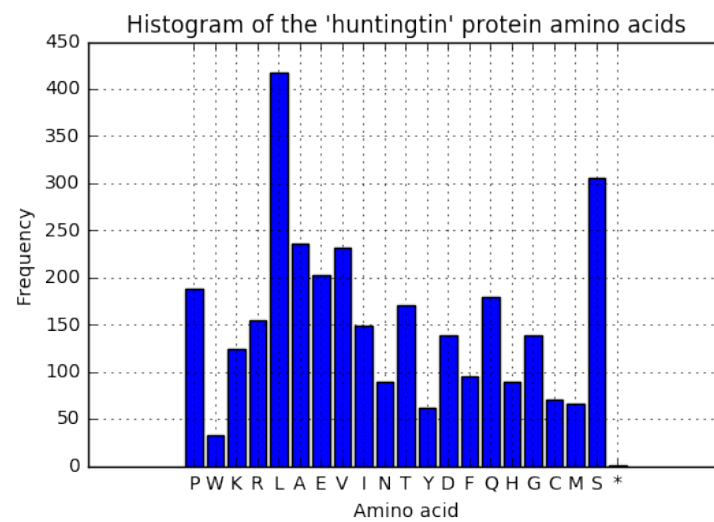


Figure 1: Histogram obtained with the Jupyter Notebook.

4 Question

Now look at the following database:

Kazusa - Codon Usage Database

The *Codon Usage Database* lists the frequency which each codon is used in a species (different species prefer different codons). Sequences which have too many rarer codons result in slowing down transcription and inhibition of protein expression - in extreme cases, rare codons are thought to introduce transcription errors when the rare tRNA is not available.

If you were to try and express your human cDNA sequence in yeast (*Saccharomyces cerevisiae*), which codons in your sequence might cause problems for expression?

Note there is no hard threshold, but generally codons with 1% usage or less are considered rare.

According to the [Saccharomyces cerevisiae DB page](#), yeast have **6534504** codons. Below is a table with the frequency of each codon in the format *[triplet] [frequency: per thousand] ([number])*.

UUU 26.1 (170666) UUC 18.4 (120510) UUA 26.2 (170884) UUG 27.2 (177573)	UCU 23.5 (153557) UCC 14.2 (92923) UCA 18.7 (122028) UCG 8.6 (55951)	UAU 18.8 (122728) UAC 14.8 (96596) UAA 1.1 (6913) UAG 0.5 (3312)	UGU 8.1 (52903) UGC 4.8 (31095) UGA 0.7 (4447) UGG 10.4 (67789)
CUU 12.3 (80076) CUC 5.4 (35545) CUA 13.4 (87619) CUG 10.5 (68494)	CCU 13.5 (88263) CCC 6.8 (44309) CCA 18.3 (119641) CCG 5.3 (34597)	CAU 13.6 (89007) CAC 7.8 (50785) CAA 27.3 (178251) CAG 12.1 (79121)	CGU 6.4 (41791) CGC 2.6 (16993) CGA 3.0 (19562) CGG 1.7 (11351)
AUU 30.1 (196893) AUC 17.2 (112176) AUA 17.8 (116254) AUG 20.9 (136805)	ACU 20.3 (132522) ACC 12.7 (83207) ACA 17.8 (116084) ACG 8.0 (52045)	AAU 35.7 (233124) AAC 24.8 (162199) AAA 41.9 (273618) AAG 30.8 (201361)	AGU 14.2 (92466) AGC 9.8 (63726) AGA 21.3 (139081) AGG 9.2 (60289)
GUU 22.1 (144243) GUC 11.8 (76947) GUA 11.8 (76927) GUG 10.8 (70337)	GCU 21.2 (138358) GCC 12.6 (82357) GCA 16.2 (105910) GCG 6.2 (40358)	GAU 37.6 (245641) GAC 20.2 (132048) GAA 45.6 (297944) GAG 19.2 (125717)	GGU 23.9 (156109) GGC 9.8 (63903) GGA 10.9 (71216) GGG 6.0 (39359)

By examining the table, one can see that **UAA** (1.1%), **UAG** (0.5%), and **UGA** (0.7%) are somewhat rare triplets (for yeast). All these three codons correspond to the **termination** codon - see [RNA codon table](#).

This is not a problem however, because despite the low percentage of these codons, there are thousands of each available, which by far cover the needs of the protein expression.

Another codon which might be considered *rare* is **CGG**, but again that should not be a problem because the expression does not demand a number of that triplet greater than its frequency in yeast.

Having said that, suppose the frequency of **CGG** was much lower - so low that it could introduce transcription errors. One possible work around is to reverse engineer the *mRNA*. **CGG** codes *Arginine*, which is also coded by **CGU**, **CGC**, and **CGA** - which have a much higher frequency (check table above). One could replace **CGG** occurrences in the *mRNA* with one of those other three *Arginine* encoding codons, and that should solve the transcription errors.

5 Question

We now turn to sequence comparison and alignment. You are given the following coding sequence fragments. They encode a homologous proteins in different species, sequence 2 is human (1 extra mark if you can give the gene name and the most likely species for the other sequences). The sequences are aligned to the correct reading frame:

1. CTGAAGCGGGAGGCTGAGACGCTGCGGGAGCGGGAGGGC
2. CTCAAGCGTGAGGCCGAGACCCTACGGGAGCGGGAAGGC
3. GAAGAGCTGAAGAGAGAGGCTGACAATTTAAAGGACAGA
4. AACGAGGAGCTCAAGCGAGAAGCTGATACGCTGAAGGAC

- 5.1 Sequences 1 and 2 differ slightly. How does the resulting protein differ? Could this have functional implications?**

TODO

- 5.2 Now use the Needleman Wunsch algorithm to compare sequence 1 to sequences 3 and 4. Use the scoring: match +2, mismatch -1, indel -1. Perform at least one of these on paper (or both if you wish). On paper, use the first three codons only.**

TODO

- 5.3 Comparing the bare sequences, what can you conclude about the relatedness of the species?**

TODO

5.4 Extra mark

TODO

Sequence	Gene name	Species
1	?	?
2	?	human
3	?	?
4	?	?

6 Sources

Genetics Home Reference

<https://ghr.nlm.nih.gov/condition/huntington-disease>

Wikipedia

https://en.wikipedia.org/wiki/Huntington%27s_disease