# Hand-in 3

Michael Iversen
Student ID: 201505099

December 5, 2022

    In this hand-in, I implement a hidden Markov model following the Codon model from the lectures. The model contains 5 hidden states: Non-coding, start codon, coding, stop codon, reverse start codon, reverse coding and reverse stop codon. The non-coding state emits from the alphabet $\mathcal{A} = \{A, G, C, T\}$ while all other hidden states emits from the alphabet $\mathcal{A} \times \mathcal{A} \times \mathcal{A} = \{AAA, AAG, AAC, AAT, AGA, ..., TTT\}$. The final model is illustrated in Fig. 1 along with the non-zero transition and emission probabilities. The transition and emission probabilities which are zero have been omitted in this figure. I obtain the model from training-by-counting without hard-coding any transition or emission probabilities. A priori, every transition between hidden states and emission of any symbols are allowed. However, after training, some transitions and emission have probability zero because they didn't appear during training.

    I translate the given annotations of non-coding $N$, coding $C$ and reverse coding $R$ into a list of hidden states. I represent the non-coding state as "0", start codon as "1", coding state as "2", stop codon as "3", reverse start codon as "4", reverse coding state as "5" and reverse stop codon as "6". The non-coding part of the annotations are directly translated $N \rightarrow 0$. The start codons are identified as a shift from $N$ to $C$. Similarly, the stop codons are identified as a shift from $C$ to $N$. The coding state is determined as part of the annotations where three letters are $C$ with both the preceding 3 and following 3 letters are all $C$. The reverse start codons, reverse stop codons and reverse coding states are identified using similar considerations. The code snippet below illustrates this translation.

```python
def str_to_list(z: str) -> list[int, ...]:
    z_list = []
    n = 0
    while n < len(z):
        if z[n] == 'N':
            z_list.append(0)
            n += 1
        elif z[n] == 'C':
            if 'N' in z[n - 3: n]:
                z_list.append(1)
            elif 'N' in z[n + 1: n + 4]:
                z_list.append(3)
            else:
                z_list.append(2)
            n += 3
        elif z[n] == 'R':
            if 'N' in z[n - 3: n]:
                z_list.append(4)
            elif 'N' in z[n + 1: n + 4]:
```
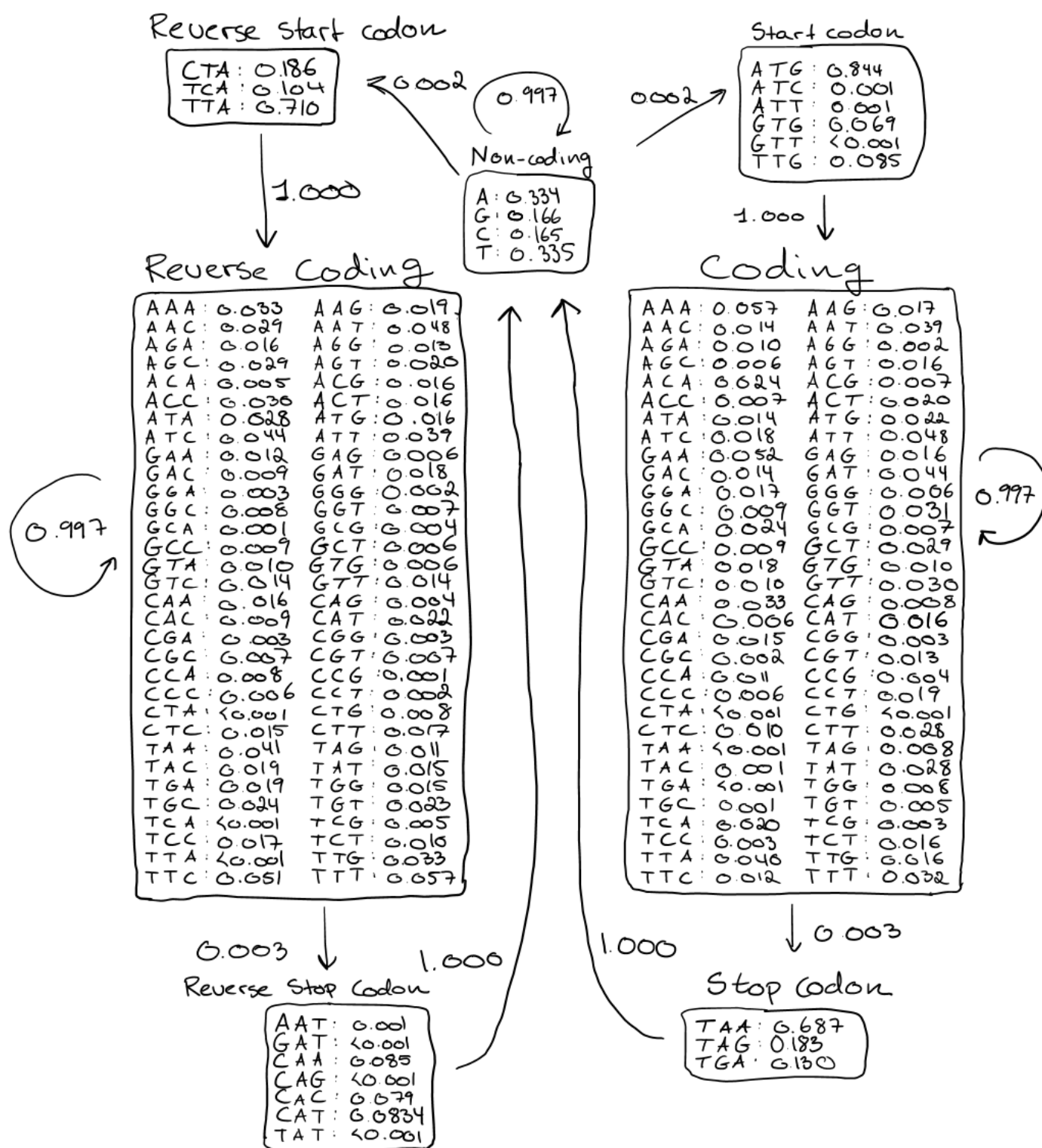
**Reverse start codon**

CTA : 0.186
TCA : 0.104
TTA : 0.710

< 0.002

0.997

0.002 ↗

**Start codon**

ATG : 0.844
ATC : 0.001
ATT : 0.001
GTG : 0.069
GTT : < 0.001
TTG : 0.085

**Non-coding**

A : 0.334
G : 0.166
C : 0.165
T : 0.335

1.000

1.000 ↓

**Reverse Coding**

| | | | |
|---|---|---|---|
| AAA : 0.033 | AAG : 0.019 |
| AAC : 0.029 | AAT : 0.048 |
| AGA : 0.016 | AGG : 0.013 |
| AGC : 0.029 | AGT : 0.020 |
| ACA : 0.005 | ACG : 0.016 |
| ACC : 0.030 | ACT : 0.016 |
| ATA : 0.028 | ATG : 0.016 |
| ATC : 0.044 | ATT : 0.039 |
| GAA : 0.012 | GAG : 0.006 |
| GAC : 0.009 | GAT : 0.018 |
| GGA : 0.003 | GGG : 0.002 |
| GGC : 0.008 | GGT : 0.007 |
| GCA : 0.001 | GCG : 0.004 |
| GCC : 0.009 | GCT : 0.006 |
| GTA : 0.010 | GTG : 0.006 |
| GTC : 0.014 | GTT : 0.014 |
| CAA : 0.016 | CAG : 0.004 |
| CAC : 0.009 | CAT : 0.022 |
| CGA : 0.003 | CGG : 0.003 |
| CGC : 0.007 | CGT : 0.007 |
| CCA : 0.008 | CCG : 0.001 |
| CCC : 0.006 | CCT : 0.002 |
| CTA : < 0.001 | CTG : 0.008 |
| CTC : 0.015 | CTT : 0.017 |
| TAA : 0.041 | TAG : 0.011 |
| TAC : 0.019 | TAT : 0.015 |
| TGA : 0.019 | TGG : 0.015 |
| TGC : 0.024 | TGT : 0.023 |
| TCA : < 0.001 | TCG : 0.005 |
| TCC : 0.017 | TCT : 0.016 |
| TTA : < 0.001 | TTG : 0.033 |
| TTC : 0.051 | TTT : 0.057 |

0.997

**Coding**

| | | | |
|---|---|---|---|
| AAA : 0.057 | AAG : 0.017 |
| AAC : 0.014 | AAT : 0.039 |
| AGA : 0.010 | AGG : 0.002 |
| AGC : 0.006 | AGT : 0.016 |
| ACA : 0.024 | ACG : 0.007 |
| ACC : 0.007 | ACT : 0.020 |
| ATA : 0.014 | ATG : 0.022 |
| ATC : 0.018 | ATT : 0.048 |
| GAA : 0.052 | GAG : 0.016 |
| GAC : 0.014 | GAT : 0.044 |
| GGA : 0.017 | GGG : 0.006 |
| GGC : 0.009 | GGT : 0.031 |
| GCA : 0.024 | GCG : 0.007 |
| GCC : 0.009 | GCT : 0.029 |
| GTA : 0.018 | GTG : 0.010 |
| GTC : 0.010 | GTT : 0.030 |
| CAA : 0.033 | CAG : 0.008 |
| CAC : 0.006 | CAT : 0.016 |
| CGA : 0.015 | CGG : 0.003 |
| CGC : 0.002 | CGT : 0.013 |
| CCA : 0.011 | CCG : 0.004 |
| CCC : 0.006 | CCT : 0.019 |
| CTA : < 0.001 | CTG : < 0.001 |
| CTC : 0.010 | CTT : 0.028 |
| TAA : < 0.001 | TAG : 0.008 |
| TAC : 0.001 | TAT : 0.028 |
| TGA : < 0.001 | TGG : 0.008 |
| TGC : 0.001 | TGT : 0.005 |
| TCA : 0.020 | TCG : 0.003 |
| TCC : 0.003 | TCT : 0.016 |
| TTA : 0.040 | TTG : 0.016 |
| TTC : 0.012 | TTT : 0.032 |

0.997

0.003 ↓    1.000    1.000

0.003 ↓

**Reverse Stop Codon**

AAT : 0.001
GAT : < 0.001
CAA : 0.085
CAG : < 0.001
CAC : 0.079
CAT : 0.0834
TAT : < 0.001

**Stop Codon**

TAA : 0.687
TAG : 0.183
TGA : 0.130

Figure 1: Model of hidden Markov model with transition and emission probabilities. Only non-zero probabilities are shown.

| Known genomes | Only Cs | Only Rs | Both |
|---|---|---|---|
| **Genome 1** | 0.7144 | 0.7675 | 0.6246 |
| **Genome 2** | 0.7549 | 0.7411 | 0.6303 |
| **Genome 3** | 0.7870 | 0.7751 | 0.6760 |
| **Genome 4** | 0.7642 | 0.7379 | 0.6248 |
| **Genome 5** | 0.7930 | 0.7390 | 0.6478 |

| Unknown genomes | Only Cs | Only Rs | Both |
|---|---|---|---|
| **Genome 6** | 0.7439 | 0.7513 | 0.6356 |
| **Genome 7** | 0.7630 | 0.7384 | 0.6380 |
| **Genome 8** | 0.7497 | 0.7049 | 0.5993 |
| **Genome 9** | 0.7209 | 0.6767 | 0.5297 |
| **Genome 10** | 0.6525 | 0.6656 | 0.4576 |

Table 1: Performance of hidden Markov model on known genomes (left) and unknown genomes (right)

```python
                z_list.append(6)
            else:
                z_list.append(5)
            n += 3
    return z_list
```

Translating a list of hidden states to an annotation in terms of $N$, $C$ and $R$ is straightforward. I simply map $0 \rightarrow N$, $1 \rightarrow CCC$, $2 \rightarrow CCC$, $3 \rightarrow CCC$, $4 \rightarrow RRR$, $5 \rightarrow RRR$ and $6 \rightarrow RRR$. This procedure is illustrated in the code snippet below.

```python
def list_to_str(z: list[int, ...]) -> str:
    def decoding(i: int) -> str:
        if i == 0:
            return 'N'
        elif i == 1 or i == 2 or i == 3:
            return 'CCC'
        else:
            return 'RRR'
    return ''.join(map(decoding, z))
```