# STRUCTURAL BIOINFORMATICS ASSIGNMENT 3

Name 1: Shruti Verma – 7022659

Name 2: Umutcan Ünaldı – 7025677

Name 3: Nazlıgül Keske - 7025902

## THEORITICAL EXERCISES

**Q1.  Using the simple Markov Model shown in Figure 1.1, write the transition probability matrix and calculate the probability that the weather for the next five days will be Rainy, Rainy, Cloudy, Sunny, Sunny, given that today is Sunny. (2 points) You can use numpy or other tools to do actual calculations (in this case you still need to provide a transition matrix and describe steps to solve this task).**
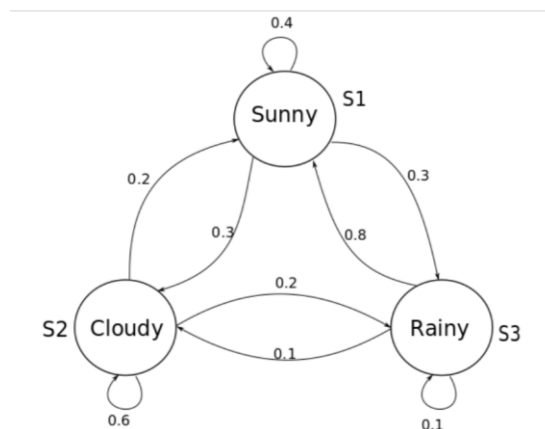


Figure 1.1: Simple Markov Model

$$
\begin{array}{cc}
& \begin{array}{ccc} S & R & C \end{array} \\
\begin{array}{c} S \\ R \\ C \end{array} &
\begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.8 & 0.1 & 0.1 \\ 0.2 & 0.2 & 0.6 \end{bmatrix}
\end{array}
$$ => Transition probability matrix, where S= Sunny, R= Rainy, C= Cloudy

## Calculating initial/ stationary state

Since the present day is 'Sunny' so we are supposing $\pi = \begin{bmatrix} S & R & C \\ 1 & 0 & 0 \end{bmatrix}$

Solving the equation $\pi A = \pi$

$$[1 \quad 0 \quad 0] \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.8 & 0.1 & 0.1 \\ 0.2 & 0.2 & 0.6 \end{bmatrix} = [0.4 \quad 0.3 \quad 0.3]$$

$$[0.4 \quad 0.3 \quad 0.3] \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.8 & 0.1 & 0.1 \\ 0.2 & 0.2 & 0.6 \end{bmatrix} = [0.46 \quad 0.21 \quad 0.33]$$

$$[0.46 \quad 0.21 \quad 0.33] \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.8 & 0.1 & 0.1 \\ 0.2 & 0.2 & 0.6 \end{bmatrix} = [0.418 \quad 0.225 \quad 0.357]$$

$$[0.418 \quad 0.225 \quad 0.357] \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.8 & 0.1 & 0.1 \\ 0.2 & 0.2 & 0.6 \end{bmatrix} = [0.4186 \quad 0.2193 \quad 0.3621]$$

$$[0.4186 \quad 0.2193 \quad 0.3621] \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.8 & 0.1 & 0.1 \\ 0.2 & 0.2 & 0.6 \end{bmatrix} = [0.4153 \quad 0.21993 \quad 0.36477]$$

$$[0.4153 \quad 0.21993 \quad 0.36477] \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.8 & 0.1 & 0.1 \\ 0.2 & 0.2 & 0.6 \end{bmatrix} = [0.415 \quad 0.219 \quad 0.365]$$

Since, **πA=π**

Therefore, $\pi = [\mathbf{0.415} \quad \mathbf{0.219} \quad \mathbf{0.365}]$

"To find the probability of the given sequence "Rainy, Rainy, Cloudy, Sunny, Sunny" **given that today is Sunny,**

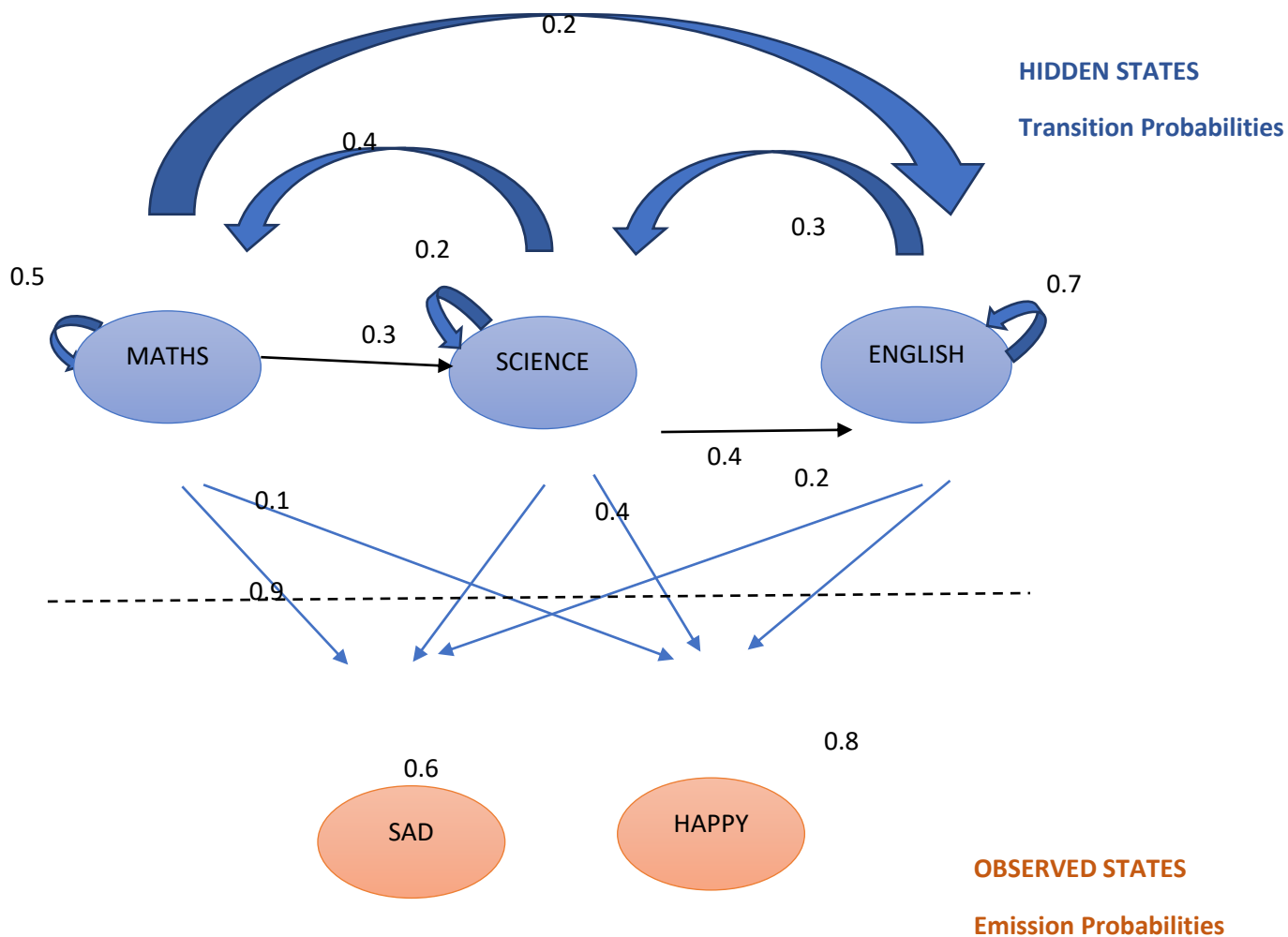⇨ P(S) x [P(R)|S] x [P(R)|R] x [P(C)|R] x [P(S)|C] x [P(S)|S]

⇨ 0.415 x 0.3 x 0.1 x 0.1 x 0.2 x 0.4

⇨ 0.0000996

Therefore, the probability of the sequence given is **0.0000996**

**Q2. Explain Hidden Markov Models (HMM) architecture with a sample diagram (don't use the one from the lecture). Make sure to mark all of the components in the diagram along with their probabilities and write down the total number of parameters.**

HIDDEN STATES

Transition Probabilities

0.2

0.4

0.2

0.3

0.5

0.7

MATHS

0.3

SCIENCE

ENGLISH

0.3

0.4

0.2

0.1

0.4

0.9

0.6

0.8

SAD

HAPPY

OBSERVED STATES

Emission Probabilities

$$
\begin{array}{ccc}
 & M & Sc & E \\
M \\
Sc & \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.4 & 0.2 & 0.4 \\ 0.0 & 0.3 & 0.7 \end{bmatrix} \\
E
\end{array}
$$

⇨ Transition Matrix

$$
\begin{array}{cc}
 & S & H \\
M \\
Sc & \begin{bmatrix} 0.9 & 0.1 \\ 0.6 & 0.4 \\ 0.2 & 0.8 \end{bmatrix} \\
E
\end{array}
$$

⇨ Emission Matrix

## PARAMETERS: -

**X** = Hidden States = Maths (M), Science (Sc), English (E)

**Y** = Observed States = Happy (H). Sad (S)

**A** = Transition probability, which is:

$$
\begin{array}{c}
\phantom{M} \\
M \\
Sc \\
E
\end{array}
\begin{array}{ccc}
M & Sc & E \\
\begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.4 & 0.2 & 0.4 \\ 0.0 & 0.3 & 0.7 \end{bmatrix}
\end{array}
$$

**B** = Output/ Emission probability, which is:

$$
\begin{array}{c}
\phantom{M} \\
M \\
S \\
E
\end{array}
\begin{array}{cc}
S & H \\
\begin{bmatrix} 0.9 & 0.1 \\ 0.6 & 0.4 \\ 0.2 & 0.8 \end{bmatrix}
\end{array}
$$

## Let's consider the example below for better explanation:-



Figure 2

**Finding the probability of the above sequence.**

*For reference:*

```
 M   Sc   E              M   Sc   E
                      M  ⎡a   d   g⎤
[A   B   C]          Sc  ⎢b   e   h⎥
                      E  ⎣c   f   i⎦
```

**Step 1:** We have to find Initial/stationary state first.

Since today the student is studying English (E), so we are supposing that $\pi = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}$

**Solving the equation $\pi A = \pi$**

$$\begin{bmatrix} 0 & 0 & 1 \end{bmatrix} \begin{matrix} M & Sc & E \\ \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.4 & 0.2 & 0.4 \\ 0.0 & 0.3 & 0.7 \end{bmatrix} \end{matrix} = \begin{bmatrix} 0 & 0.3 & 0.7 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0.3 & 0.7 \end{bmatrix} \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.4 & 0.2 & 0.4 \\ 0.0 & 0.3 & 0.7 \end{bmatrix} = \begin{bmatrix} 0.12 & 0.27 & 0.61 \end{bmatrix}$$

$$\begin{bmatrix} 0.12 & 0.27 & 0.61 \end{bmatrix} \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.4 & 0.2 & 0.4 \\ 0.0 & 0.3 & 0.7 \end{bmatrix} = \begin{bmatrix} 0.168 & 0.273 & 0.559 \end{bmatrix}$$

$$\begin{bmatrix} 0.168 & 0.273 & 0.559 \end{bmatrix} \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.4 & 0.2 & 0.4 \\ 0.0 & 0.3 & 0.7 \end{bmatrix} = \begin{bmatrix} 0.1932 & 0.2727 & 0.5341 \end{bmatrix}$$

$$\begin{bmatrix} 0.1932 & 0.2727 & 0.5341 \end{bmatrix} \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.4 & 0.2 & 0.4 \\ 0.0 & 0.3 & 0.7 \end{bmatrix} = \begin{bmatrix} 0.20568 & 0.27273 & 0.52159 \end{bmatrix}$$

$$\begin{bmatrix} 0.20568 & 0.27273 & 0.52159 \end{bmatrix} \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.4 & 0.2 & 0.4 \\ 0.0 & 0.3 & 0.7 \end{bmatrix} =$$
$$\begin{bmatrix} 0.211932 & 0.272727 & 0.51534 \end{bmatrix}$$

$$\begin{bmatrix} 0.211932 & 0.272727 & 0.51534 \end{bmatrix} \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 0.4 & 0.2 & 0.4 \\ 0.0 & 0.3 & 0.7 \end{bmatrix} = \begin{bmatrix} 0.218 & 0.273 & 0.51 \end{bmatrix}$$

---

Since, $\pi A = \pi$

Therefore, $\pi = \begin{bmatrix} 0.218 & 0.273 & 0.51 \end{bmatrix}$

---

**P(Y = H, H, S ; X = E, Sc, E)**

$P(X_1 = E) \times P(Y_1 = H \mid X_1 = E) \times P(Y_2 = H \mid X_2 = Sc) \times P(Y_3 = S \mid X_3 = E) \times P(X_2 = Sc \mid X_1 = E) \times P(X_3 = E \mid X_2 = Sc)$

⇨ 0.51× 0.8 x 0.4 x 0.2 x 0.3 x 0.4

⇨ 0.00391

Therefore, probability for the above sequence is **0.00391**

**Q3. Explain sequence logos and compare them to sequence profiles. Where can they be used? (1 point)**

Sequence logos are representation of conserved sequence along a series of multiple sequence alignment of genetically similar sequences in a single graphical form. It is created by a consensus sequence of all the given sequences. The higher the frequency of occurrence of a certain base pair in a given sequence alignment of DNA, RNA or protein, higher it will be the height of that particular base in the graphical representation. This way we can easily see which bases or amino acids are more common compared to others in an alignment.

On the other hand, sequence profiles give much more information about consensus sequences. This is due to information coming from occurrences. In addition to consensus sequence, sequence profiles also calculate the occurrence of amino acids or nucleotides for each of the position as exact numbers. This gives an estimate prediction of probability of which amino acid or nucleotide can occur on that position.

Both methods can be used any of the multiple sequence alignments. They both can make it easier for users to detect the abundance of gene or amino acids on particular positions.

**Q4. What are the differences between the traditional artificial neural networks and deep learning? (1 point)**

| | Traditional Artificial neural network | Deep learning |
|---|---|---|
| 1. | An artificial neural network imitates the actual human brain neurons and trains itself to give best possible output according to the real world trained based on assorted sets of algorithms designed to mimic how neurons perceive and react to sensory data | Deep learning imitates the data processing techniques of human brain just like ANN but it is rather trained based on representation of data which is unstructured unlike ANN. |
| 2. | A neural network basically consists of 3 layers usually with 1 hidden layer along with 1 input layer and 1 output layer. | A Deep learning system can be referred to as an ANN with more than 3 layers including1 input layer and 1 output layer |
| 3. | Traditional artificial neural networks are a subset of the field of Machine learning. | Deep learning is a subset of the field of Artificial neural networks in the area of Machine learning. |
| 4. | It is applicable to detect objects and facial character recognition. Identification and classification of text recognition according to the relevant certain categories. | It can be in speech and textual recognition and categorization and for generation of HD videos depending om the observations on low quality image and footages which can further be used for development of high-quality images of historic data. It can also be sued for digital marketing by showing advertisements based on the surfer's previous search history. |

**PROGRAMMING EXERCISES**

**1**- 1A9U.pdb is used for this assignment.

10 proteins with 40 to 70 percent identity is chosen for multiple alignment.



**2-** The file is downloaded as SeqDump.txt and FASTA of original .pdb is added on top of it.

Input is selected for multiple sequence alignment.



Result of alignment.



**3-**

```
# Panda is called for DataFrame
import pandas as pd
# Bio and Bio.Align is specifically used for consensus calculations
from Bio import AlignIO as al
from Bio.Align import AlignInfo as alI


# Alignment is called with the format "clustal"
alignment = al.read("clustalo-E20211130-094437-0140-87358764-
p2m.clustal_num", "clustal")
```

```
# Summary info is assigned
align = alI.SummaryInfo(alignment)
# Dumb_consensus function of the AlignInfo module is used to generate
consensus
consensus = align.dumb_consensus()
print("Consensus of the clustal is:" + "\n" + consensus)
```

Gives the output as:

Consensus of the clustal is:

MGSSHHHHHHSSGLVPRGXSXXXXXRSGFYRQEVTKTAWEVRAVYRDLXPVGSGA
YGAVCSAVDGRTGAKVAIKKLYRPFQSELFAKRAYRELRLLKHMRHENVIGLLDVF
TPDETLDDFXDFYLVMPFMGTDLGKLMKHEKLGEDRIQFLVYQMLKGLRYIHXAGI
IHRDLKPGNLAVNEDCELKILDFGLARQADSEMTGYVVTRWYRAPEVILNWMXYT
QTVDIWSVGCIMAEMITGKTLFKGSDHLDQLKEIMKVTGTPPAEFVQRLQSDEAKN
YMKGLPELEKKDFASILTNASPLAVNLLEKMLVLDAXXRXTAGEALAHPYFESLHD
TEDEPQAVQKYDDSFDXXDRTLDEWKRVTYKEVLSFKPPRQLGARVSKETPL