

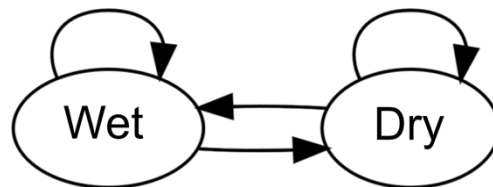
Lab 7: Limiting Behavior of Markov Chains

Instructions: Read the prompts and use MATLAB to answer the questions. Submit your solution as a single script. A template script is available on Canvas.

Part 1. Simple Weather System

Weather is a complex physical phenomenon. Nevertheless, some tendencies of weather patterns can be captured by relatively simple models. For example, in 1964, Leonard L. Weiss developed a simple Markov chain probability model of wet and dry days¹. The model was very effective in predicting the duration of wet and dry sequences across a diverse set of locations. Let's develop a similar Markov chain model of the weather in Davis and use it as an example to learn about the *limiting distribution* of Markov chains.

Consider a simple Markov chain model of the weather with two states: dry and wet.



A student in Davis, California recorded the daily weather in 2018. The first day of record was dry; after that, 270 days were observed to be dry the day after a dry day; 36 days were observed to be dry the day after a wet day; 28 days were observed to be wet the day after a dry day; and 30 days were observed to be wet the day after a wet day (in reality, we know that Davis experiences many other types of weather – cloudy, foggy, etc. – but for this example, any non-rainy day is called “dry” for simplicity).

With this data in hand, we can construct a matrix of counts for the transitions between the two states in our system:

$$\begin{bmatrix} \text{\#of dry days after a dry day} & \text{\#of wet days after a dry day} \\ \text{\#of dry days after a wet day} & \text{\#of wet days after a wet day} \end{bmatrix} = \begin{bmatrix} 270 & 30 \\ 28 & 30 \end{bmatrix}$$

This matrix informs us how many times each feasible weather transition occurred during 2018, but it's not a transition matrix yet. This is because a transition matrix holds the probability of transitioning from one state to another, not the total counts from a series of observations. So, we need to convert our observed counts into probabilities.

¹[https://doi.org/10.1175/1520-0493\(1964\)092<0169:SOWODD>2.3.CO;2](https://doi.org/10.1175/1520-0493(1964)092<0169:SOWODD>2.3.CO;2)

To do this, we need to divide each count by the number of times the system was in the previous state. Let's look at an example – the top left element of the transition matrix should give us the probability of the weather staying dry after a dry day. We know that 270 days were dry after a dry day, and 30 days were wet after a dry day. Thus, of the 300 days which were dry, $270/300=0.9$ remained dry, and $30/300=0.1$ transitioned to rain. These values are the transition probabilities in the first row of our transition matrix! Applying the same procedure to wet days, we arrive at:

$$P = \begin{array}{cc} & \begin{array}{c} \text{Second Day} \\ \text{Dry} \quad \text{Wet} \end{array} \\ \begin{array}{c} \text{First} \\ \text{Day} \end{array} \begin{array}{c} \text{Dry} \\ \text{Wet} \end{array} & \begin{bmatrix} 0.90 & 0.10 \\ 0.48 & 0.52 \end{bmatrix} \end{array}$$

Question 1: Encode this transition matrix in MATLAB and save it into a variable named `P`.

Question 2: Recall that the distribution of states at a point in time can be computed using the **Chapman-Kolmogorov equation**, $\Pi(t) = \Pi(0)P^t$. Compute a large power of our transition matrix, P (a power larger than 100). What does the resulting matrix look like?

Question 3: Assume a weather sequence starts as a dry day, meaning that $\Pi(0) = [1, 0]$. To what distribution does the system converge to as t gets large? If a weather sequence starts as a wet day instead, is the distribution the same?

Taking the transition matrix to a large power via the Chapman-Kolmogorov equation is only guaranteed to converge to a unique limiting distribution if the Markov chain is irreducible, aperiodic, and positively recurrent. To arrive at stationary solutions more generally, we can use a different approach leveraging some linear algebra. Remember that a limiting distribution is one which remains constant over time. In other words, if Π_L is a limiting distribution, it must be an eigenvector of the matrix P , which is defined as:

$$\Pi_L \lambda = \Pi_L P$$

where λ is a scalar constant, referred to as the corresponding eigenvalue of the eigenvector. In the case of a transition matrix like the one characterizing the weather in Davis, the eigenvector associated with an eigenvalue of $\lambda = 1$ is the limiting distribution ($\Pi_L = \Pi_L P$).

MATLAB is well-optimized to compute eigenvectors and eigenvalues of a matrix. The `eig` command is what will perform the analysis for you. If you want to compute the eigenvalues and eigenvectors of some matrix `A`, use the command `[V,D]=eig(A)`. The output `V` will be a matrix where each column is an eigenvector. The output `D` will be a matrix with the eigenvalues on the diagonal.

Question 4: Use the `eig` function in MATLAB to compute the eigenvalues and eigenvectors of the transition matrix describing weather in Davis. (You'll need to use the transpose of `P` when you call the function in order to orient the matrix properly for MATLAB, which by default returns *right* eigenvectors instead of *left* eigenvectors as written above).

Question 5: The eigenvector associated with the limiting distribution of our Markov chain will be held in the first column of the eigenvector matrix, \mathbf{v} (in other words, the first eigenvector is $\mathbf{v}(:, 1)$). What is this eigenvector?

Question 6: A limiting distribution, by definition, must sum to 1. Does the eigenvector from Question 5 sum to 1?

Question 7: A vector can be scaled to sum to 1 by dividing it by its sum. Divide $\mathbf{v}(:, 1)$ by its sum and save the result in a new vector – what result do you get? Compare this result to your answers from Question 3.

Question 8: The limiting distribution of the Markov chain tells us what proportion of days are expected to be dry and wet, respectively. Do you think these proportions would accurately capture the weather trends in Davis over the course of a full year? How might the seasonality of weather in Davis affect the quality of the model? How might the model be improved?

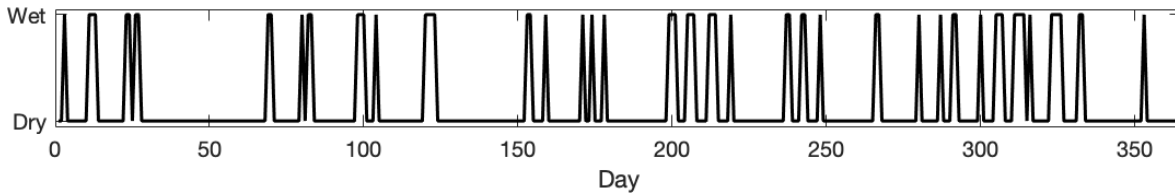
A couple lines of code in MATLAB can simulate sequences according to our Markov model of weather. Starting with an initial state, we just need to sample each next state according to the transition probabilities of the state before it. Setting the initial state to be a dry day, we could produce sequences with `rand` and a couple `if`-statements. For example, to simulate 365 days of weather from our model:

```
P = [0.9, 0.10; 0.48, 0.52]; % Transitions
T = 365;
X = NaN(1,T);
X(1) = 1;
for i = 2:T
    probs = P(X(i-1),:);
    X(i) = binornd(1, probs(2))+1;
end
display(X)
```

Which gives a sequence:

1 1 2 1 1 1 1 1 1 1 2 2 ...

Plotting the entire sequence of wet and dry days with `plot` gives a figure that looks something like this:



Question 9: Simulate a sequence of 365 days with the given model of Davis weather.

Suppose that we wanted to quantify how the weather *converges* to the limiting distribution of states. Given a sequence of states \mathbf{x} , we can compute the *proportion of wet days* over time with `cumsum(X-1) ./ (1:T)`.

Question 10: Compute the proportion of wet days over time for your simulated sequence. Plot the proportion as it changes over time with `plot` – does the proportion of wet days converge to the prediction from the limiting distribution?

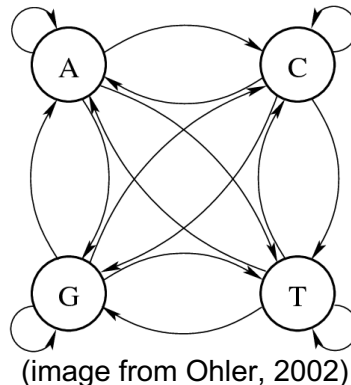
Question 11: Consider a different location called "Random Town" where the weather transitions are characterized by a Markov chain with $P = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$.

11a: What is the limiting distribution of this model?

11b: Simulate a sequence from this model and plot its proportion of wet days over time (as in Question 10). Does this system seem to converge faster or slower than the Davis weather system?

Part 2. DNA Sequence of *E.coli* genome

A slightly more sophisticated application of a Markov chain model is to a sequence of DNA. DNA is comprised by four nucleotides (A, C, T, G), and as such, a Markov chain model can naturally be formulated using these four states. Each state can transition to any of the four states, so the graphical representation of the model might look something like this:



In order to characterize a DNA sequence using a Markov chain model, one could count how many times each of the feasible state transitions occur in the sequence. This is analogous to counting the transitions in weather from our previous example; however, a Markov chain model of DNA would have 16 possible transitions (4 nucleotides with 4 transitions each) instead of the four transitions that arise from the two-state dry-wet model. These 16 transitions are referred to as “dinucleotide pairs.”

For instance, a researcher may wish to characterize the DNA genome of the bacteria *E. coli*. This bacterium has been studied extensively since the late 1890s when it was initially discovered by German-Austrian pediatrician Theodor Escherich. The complete genome of *E. coli* was only sequenced relatively recently, however, in the 1990s. Nowadays, it’s easy to download the genome from the Internet. The complete genome of a common *E. coli* strain (K12) is available with this lab as a `.fasta` file.

An interesting property of the *E. coli* genome is how uniformly the four nucleotides are utilized in the sequence. If you compute the proportion of the genome that is comprised by A’s, by T’s, by G’s, and by C’s, you’ll find that each nucleotide is responsible for almost exactly 25% of the sequence. What does this uniformity say about nucleotide transitions, if anything? Is the same uniformity present in dinucleotide pairs?

Let’s investigate. The code below (1) reads in the *E. coli* genome found in `Ecoli-k12-genome.fasta`, (2) counts the number of times each dinucleotide pair occurs, and (3) normalizes the dinucleotide counts to produce a transition matrix, based on the DNA sequence. Note that we have mapped the letter “A” to the first row/column of the transition matrix, “T” to the second row/column, “G” to the third row/column, and “C” to the fourth.

```

file = fastaread('Ecoli-k12-genome.fasta');
seq = file.Sequence;

di_nuc_counts = zeros(4);

for i = 1:length(seq)-1

    di_nuc = seq(i:i+1); % Read dinucleotide

    row = nuc_to_index(di_nuc(1)); % Convert letters to index for matrix
    col = nuc_to_index(di_nuc(2)); % Convert letters to index for matrix

    di_nuc_counts(row, col) = di_nuc_counts(row, col) + 1;

end

transition_matrix = di_nuc_counts./sum(di_nuc_counts, 2);
display(transition_matrix);

function index = nuc_to_index(nuc)

switch nuc
    case {'A', 'a'}
        index = 1;
        return
    case {'T', 't'}
        index = 2;
        return
    case {'G', 'g'}
        index = 3;
        return
    case {'C', 'c'}
        index = 4;
        return
end
end

```

Question 12: Download the *E. coli* genome file and run the code above to obtain the transition matrix. Are the dinucleotide pairs used uniformly? If not, which dinucleotide pairs are used the most?

Question 13: Find the limiting distribution of this Markov chain model. Do you observe the nucleotides are used about equally?

Assume that a researcher is placed in a certain predicament. The researcher has two DNA sequences, but they've been mis-labeled, and they are unsure which sequence belongs to which organism. Initially, they characterize the sequences by their nucleotide composition to distinguish them, but it is determined that both sequences have the same composition, which is 14.3% A, 14.3% T, 35.7% G, and 35.7% C.

Having been unable to distinguish the sequences by their composition, the researcher instead decides to check for differences in their nucleotide transitions. One of the organisms, *E. atrepeatum*, is known to contain some regions where many A's appear one after another and other regions where many T's appear one after another. The other organism, *E. norepeatium*, does not have the same propensity to have A's and T's appear in such long stretches, respectively. The researcher believes this difference may assist in distinguishing the genomes.

The two mis-labeled genomes were posted with this lab as "**Organism_A.fasta**" and "**Organism_B.fasta**." Download their genomes and answer the questions below.

Question 14: Use the code provided for the *E. coli* example to analyze the nucleotide transitions that comprise these two genomes. What does the transition matrix look like for each genome?

Question 15: For each genome's transition matrix, confirm that the overall composition tends to the values given in the text above by computing their limiting distribution.

Question 16: Does analyzing the dinucleotide counts for these genomes allow us to distinguish which genome belongs to which organism? Which genome (A or B) do you believe belongs to *E. atrepeatum*?