

MSiA-400 Everything Starts with Data

Lab Assignment 1

Due date: Thursday November 1, 12 pm

EXERCISE INSTRUCTIONS: Please submit one report file that includes: short answer, related code and print for each problem if necessary. Push your answers to Github (required) and Canvas (optional).

Problem 1

In *Markov100.txt*, the one step transition probability matrix for a Markov chain with 100 states (State 1 to State 100) is given. Note that the data has no heading.

Name of the data set	Markov100
Number of rows	100
Number of columns	100

Problem 1(a)

Suppose we are at State 1 now. Find and display the probability of being in State 5 after 10 transitions.

Problem 1(b)

Suppose we are at one of States 1, 2, and 3 with equal probabilities. Find and display the probability of being in State 10 after 10 transitions.

Problem 1(c)

Find the steady state probability of being in State 1.

Problem 1(d)

Find the mean first passage time from State 1 to State 100.

Problem 2

You are asked to analyze the data from an website with 8 pages (Page 1 - Page 8). Let us assume that there is a virtual page (Page 9) that a visitor must automatically visit when the visitor leaves the website. The visitors always start their visit from Page 1. Let us formulate a Markov chain for this website. The states are defined as

$$S_i = \text{visitor is at Page } i, i = 1, \dots, 9.$$

For example, suppose that a visitor enters the website (hence visit Page 1), moves to Page 3, Page 5, and then leave the website, sequentially. Then, the user visits States S_1, S_3, S_5 , and S_9 , sequentially.

Please find the attached data *webtraffic.txt*. The data includes the record of 1000 visitors (rows). The data has 81 columns labeled as $t_{11}, t_{12}, \dots, t_{19}, t_{21}, t_{22}, \dots, t_{29}, \dots, t_{91}, t_{92}, \dots, t_{99}$. The label t_{ij} represents the transition from State i to State j , for $i = 1, \dots, 9$ and $j = 1, \dots, 9$. For example, t_{12} is the transition from State 1 to State 2, and t_{84} is the transition from State 8 to State 4. For each visitor (row), it has 1 for column t_{ij} if the visitor makes transition from State i to State j , and it has 0 elsewhere. For example, if a visitor

visits States S_1, S_3, S_5 , and S_9 , sequentially, then the corresponding row has 1 for columns t_{13}, t_{35}, t_{59} and 0 elsewhere.

The summary of the data set is below.

Name of the data set	webtraffic
Type of data	binaries (0,1)
Number of rows	1000
Number of columns	81

Problem 2(a)

Construct 9 by 9 matrix **Traffic** that counts total traffic between State i to State j for all $i = 1, \dots, 9$ and $j = 1, \dots, 9$. Display **Traffic**.

Hint `colSums()` adds all rows for each column.

Problem 2(b)

Observe that **Traffic** has 0's in row 9 and 0's in column 1. Set **Traffic**[9,1]=1000. Construct the one step transition probability matrix **P** and display it.

Problem 2(c)

Calculate and display the steady state probability vector **Pi**.

Problem 2(d)

The following table presents the average time that the visitors spend on each page.

Page	1	2	3	4	5	6	7	8
Avg(minute)	0.1	2	3	5	5	3	3	2

Calculate and display the average time a visitor spend on the website (until she leaves).

Problem 2(e)

In the output of Problem 2(c), observe that Pages 3 and 4 are one of the most crowded pages except Pages 1 and 9. To balance the traffic, the owner of the website decided to create links from Page 2 to Pages 6,7 (hence, from State 2 to States 6,7). By adding the links, the owner anticipates that, from Page 2, 30% of the current outgoing traffic to State 3 would move to State 6, and 20% of the current outgoing traffic to State 4 would move to State 7. Calculate new steady state probability vector **Pi2** to check the effect of the new links. Decide if the link helped balancing the traffic by comparing the variance of **Pi** and **Pi2**.

Hint Start with matrix **Traffic** from Problem 2(a).