

1 Mutual Information in Markov Chains

If we have a Markov chain defined by the matrix \mathbf{M} (we adopt the notation in the paper and write \mathbf{M} instead of \mathbf{P}), which is *irreducible* and *aperiodic*, and has a finite state space $S = \{1, \dots, n\}$, then we have that

$$\lim_{i \rightarrow \infty} \mathbf{M}^i = \mathbf{M}_{\boldsymbol{\mu}} \quad ,$$

where $\mathbf{M}_{\boldsymbol{\mu}}$ is the matrix whose columns all consist of the unique stationary probability distribution $\boldsymbol{\mu}$. In case the reader is unfamiliar with these terms or this result, one can read them in appendix ??.

Now, let us consider two random variables X and Y , which will denote the state of the Markov chain at times t_0 and $t_0 + \tau$ respectively. We assume that we measure these variables very late in the process, where we already have that $\mathbf{M}^{t_0} \approx \mathbf{M}_{\boldsymbol{\mu}}$. We will use this fact later.

Our goal now is to quantify the mutual information $I(X; Y)$ between X and Y , that is, the discrepancy between the joint probability distribution $P_{(X,Y)}$ and the one defined by the product of the two marginalized distributions P_X and P_Y , that is $P' := P_X \otimes P_Y$. We use the Kullback-Leibler divergence, so our target expression becomes

$$I(X; Y) = D_{\text{KL}}(P_{(X,Y)} \| P_X \otimes P_Y) \quad .$$

For further details, see appendix ??.

1.1 Exponential Decay in Irreducible Aperiodic Markov Chains

Note that $I(X; Y)$ depends on the properties of M , as well as on τ . Because \mathbf{M} is irreducible and aperiodic, it follows that $|\lambda_2| < 1$. The claim is that:

Theorem 1.1 (No Power-Law in Markov Chains). *Let X and Y be random variables from an irreducible aperiodic Markov chain at times t_0 and $t_0 + \tau$ respectively. Let \mathbf{M} be the transition matrix, and let $|\lambda_2|$ denote the second largest absolute value of its eigenvalues. Then:*

$$I(X; Y) \in \mathcal{O}(|\lambda_2|^\tau) \quad .$$

Remark 1.1. Based on theorem ??, it follows that Markov chains are incapable of weak power-law behavior (and hence also strong power-law behavior).

There is a lot of math involved, so let us first get an intuition for what is going on. When considering Markov chains, we consider a set of states, say $S = \{A, B, C\}$, and for each time $t \in \mathbb{N}$ we assign a probability to the random variable $X_t \in S$. So let us consider the following Markov chain in figure 1.

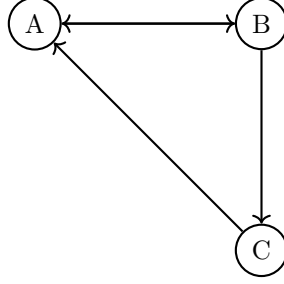


Figure 1: A simple irreducible aperiodic Markov chain. Note that if $X_{t_0} = C$, then we know that $X_{t_0+1} = A$.

If $\tau = 1$, i.e. we consider the mutual information of two consecutive states, we get a large value of $I(X, Y)$, as if X_{t_0} is either A or C , then X_{t_0+1} is uniquely determined, so we have a strong dependency between the two random variables. If, however, we have $\tau = 5$, then we can reach every state independent of the starting position. To see this, note that we can reach every state from A in four steps:

- $A \rightarrow B \rightarrow C \rightarrow A \rightarrow B$
- $A \rightarrow B \rightarrow A \rightarrow B \rightarrow C$
- $A \rightarrow B \rightarrow A \rightarrow B \rightarrow A$

The last step can then be used to go around in a cycle. If we on the other hand started at B or C , then we could go to A in one step, and consequently to every other state in the following four. Hence, the probability distribution will spread out over time and converge to the stationary probability distribution, which results in a decline of $I(X; Y)$ for increasing τ .

Because we measure our X very late in time, meaning t_0 is very large, we will have that $P(X = a) \approx \mu_a$ because of this. Similarly, we have $P(Y = b) \approx \mu_b$, since the probability distribution will only get attracted more towards μ . As we now increase τ , $P(Y = b | X = a)$ itself will converge to μ_b exactly due to the same reason. Note that $P(Y = b | X = a) = (\mathbf{M}^\tau)_{ba} \xrightarrow{\tau \rightarrow \infty} \mu_b$. And, of course, if $P(X = a, Y = b) = P(X = a) \cdot P(Y = b | X = a) = \mu_a \cdot \mu_b$, we have $I(X; Y) = 0$. Hence, in a sense the theorem describes how fast $\mathbf{M}^\tau \mathbf{p}_0$ converges to μ , or, equivalently, \mathbf{M}^τ towards \mathbf{M}_μ .

Proof. Now it's time to dive into the math. In the following, we try to reconstruct the arguments given in the paper. We also adopt the notation $P(a, b) \equiv P(X = a, Y = b)$. By definition of the Kullback-Leibler divergence, we have

$$I(X; Y) = D_{\text{KL}}(P_{(X, Y)} \| P_X \otimes P_Y) = \sum_{(a, b) \in S^2} P(a, b) \log_B \frac{P(a, b)}{P(a)P(b)} \quad .$$

The idea is now that $\log_B(\bullet)$ is *concave*. Hence, we can upper bound it by its Taylor expansion of the first degree at the point $x_0 = 1$:

$$\begin{aligned} \log_B(x) &\leq \log_B(x_0) + \log'_B(x_0)(x - x_0) \\ &= 0 + \frac{\ln'(x_0)}{\ln(B)}(x - 1) \\ &= \frac{\frac{1}{x_0}}{\ln(B)}(x - 1) \\ &= \frac{x - 1}{\ln(B)} \quad . \end{aligned}$$

For simplicity, we set $B := e$. So our expression becomes

$$\begin{aligned} I(X; Y) &\leq \frac{1}{\ln(B)} \sum_{(a, b) \in S^2} P(a, b) \left(\frac{P(a, b)}{P(a)P(b)} - 1 \right) \\ &= \sum_{(a, b) \in S^2} P(a, b) \left(\frac{P(a, b)}{P(a)P(b)} - 1 \right) \\ &= \left(\sum_{(a, b) \in S^2} P(a, b) \frac{P(a, b)}{P(a)P(b)} \right) - 1 \\ &= \left(\sum_{(a, b) \in S^2} \frac{P(a, b)^2}{P(a)P(b)} \right) - 1 \\ &=: I_R(X; Y) \quad . \end{aligned}$$

The authors of the paper coin this definition for $I_R(X; Y)$ the *rational mutual information*, as it has some useful properties. As discussed, we can approximate $P(a) \approx \boldsymbol{\mu}_a$ and $P(b) \approx \boldsymbol{\mu}_b$, and also $P(b|a) = (\mathbf{M}^\tau)_{ba}$. Thus:

$$\begin{aligned} I_R(X; Y) + 1 &= \sum_{(a, b) \in S^2} \frac{P(a, b)^2}{P(a)P(b)} \\ &= \sum_{(a, b) \in S^2} \frac{P(b|a)^2 P(a)^2}{P(a)P(b)} \\ &= \sum_{(a, b) \in S^2} \frac{\boldsymbol{\mu}_a}{\boldsymbol{\mu}_b} [(\mathbf{M}^\tau)_{ba}]^2 \quad . \end{aligned}$$

Let us now focus on $(\mathbf{M}^\tau)_{ba}$. For simplicity, we consider the case that \mathbf{M} is diagonalizable (for the general case see section 1.1.1). Note that since \mathbf{M} is irreducible and aperiodic, we have that $1 = \lambda_1 > |\lambda_2| \geq \dots \geq |\lambda_n|$. Hence, let

$$\mathbf{M} = \mathbf{B} \mathbf{D} \mathbf{B}^{-1}$$

be the diagonalization of \mathbf{M} . Of course, we immediately see that $\mathbf{M}^\tau = \mathbf{B} \mathbf{D}^\tau \mathbf{B}^{-1}$. Hence, it is easy to verify that

$$(\mathbf{M}^\tau)_{ba} = \sum_{c=1}^n \lambda_c^\tau \mathbf{B}_{bc} (\mathbf{B}^{-1})_{ca} \quad .$$

Okay, that was a lot of math. Now it is a good time to reassure ourselves what we actually have achieved. What do we expect $(\mathbf{M}^\tau)_{ba}$ to look like for $\tau \rightarrow \infty$? μ_b of course. What does \mathbf{B} look like? Well, this is very hard to tell, it at least should have a scaled version of μ in its first column. But we cannot really infer any information about \mathbf{B}^{-1} . But we know

$$\begin{aligned} \mu_b &= \lim_{\tau \rightarrow \infty} (\mathbf{M}^\tau)_{ba} \\ &= \lim_{\tau \rightarrow \infty} \sum_{c=1}^n \lambda_c^\tau \mathbf{B}_{bc} (\mathbf{B}^{-1})_{ca} \\ &= \lambda_1 \mathbf{B}_{b1} (\mathbf{B}^{-1})_{1a} \quad . \end{aligned}$$

So we know that

$$(\mathbf{M}^\tau)_{ba} = \mu_b \pm \mathcal{O}(|\lambda_2|^\tau) \quad .$$

Note that this is informal writing. It would be more precise to state that $|(\mathbf{M}^\tau)_{ba} - \mu_b| \in \mathcal{O}(|\lambda_2|^\tau)$.

This is looking promising, as this means that the discrepancy between $(\mathbf{M}^\tau)_{ba}$ and μ_b decays exponentially. The only thing left to do is translating this exponential decay to the mutual independence measure $I_R(X; Y)$. To this end, we

plug our results back into our previous equation. Thus:

$$\begin{aligned}
I_R(X; Y) &= \left(\sum_{(a,b) \in S^2} \frac{\mu_a}{\mu_b} [(\mathbf{M}^\tau)_{ba}]^2 \right) - 1 \\
&= \sum_{(a,b) \in S^2} \left(\frac{\mu_a}{\mu_b} [(\mathbf{M}^\tau)_{ba}]^2 - \mu_a \mu_b \right) \\
&= \sum_{(a,b) \in S^2} \left(\frac{\mu_a}{\mu_b} [\mu_b \pm \mathcal{O}(|\lambda_2|^\tau)]^2 - \mu_a \mu_b \right) \\
&= \sum_{(a,b) \in S^2} \left(\frac{\mu_a}{\mu_b} [\mu_b^2 \pm \mathcal{O}(|\lambda_2|^\tau)] - \mu_a \mu_b \right) \\
&= \pm \sum_{(a,b) \in S^2} \frac{\mu_a}{\mu_b} \mathcal{O}(|\lambda_2|^\tau) \quad ,
\end{aligned}$$

where we have used multiple facts about μ . For instance, $\sum_{a \in S} \mu_a = 1$ and thus $\sum_{(a,b) \in S^2} \mu_a \mu_b = 1$, as well as $0 < \mu_a < 1$ for all $a \in S$ (at least for $|S| > 1$). We now use the inequality again: We see that we can always bound $\frac{\mu_a}{\mu_b}$ from above, i.e. there exists $\alpha \in \mathbb{R}$ s.t. for all $(a,b) \in S^2$ we have $\frac{\mu_a}{\mu_b} < \alpha$. Hence:

$$\begin{aligned}
|I_R(X; Y)| &\in \sum_{(a,b) \in S^2} \frac{\mu_a}{\mu_b} \mathcal{O}(|\lambda_2|^\tau) \\
\implies |I_R(X, Y)| &\in \sum_{(a,b) \in S^2} \alpha \mathcal{O}(|\lambda_2|^\tau) \\
\implies |I_R(X, Y)| &\in n^2 \alpha \mathcal{O}(|\lambda_2|^\tau) \\
\implies |I_R(X, Y)| &\in \mathcal{O}(|\lambda_2|^\tau) \quad .
\end{aligned}$$

Of course, $I_R(X; Y) \geq 0$, so really $I_R(X; Y) \in \mathcal{O}(|\lambda_2|^\tau)$. Since $0 \leq I(X; Y) \leq I_R(X; Y)$, we also have $I(X; Y) \in \mathcal{O}(|\lambda_2|^\tau)$. \square

Remark 1.2. The above proof should also work without the approximation $P(a) \approx \mu_a$, so t_0 doesn't have to be large.

Remark 1.3. Based on the proof, we see that if the distance between \mathbf{M}^τ and \mathbf{M}_μ experiences exponential decay, we can translate this exponential decay to the mutual information measure $I_R(X, Y)$. Note that we have already established that *all* irreducible aperiodic Markov chains have this property in remark ??.

1.1.1 The Defective Case

Nonetheless, we will prove the case that \mathbf{M} is not diagonalizable separately and establish the connection to λ_2 . The idea is that while not every matrix is diagonalizable, every square matrix over the complex numbers can be put into *Jordan normal form*, which resembles diagonalization. In this form, the matrix is nearly diagonal, except that for each repeated eigenvalue, there may be 1s on the superdiagonal (just above the main diagonal), indicating the presence of generalized eigenvectors.

For example, if there are only three distinct eigenvalues and λ_2 is threefold degenerate, the the Jordan form of \mathbf{M} would be

$$\mathbf{B}^{-1}\mathbf{M}\mathbf{B} = \begin{bmatrix} \lambda_1 & 0 & 0 & 0 & 0 \\ 0 & \lambda_2 & 1 & 0 & 0 \\ 0 & 0 & \lambda_2 & 1 & 0 \\ 0 & 0 & 0 & \lambda_2 & 0 \\ 0 & 0 & 0 & 0 & \lambda_3 \end{bmatrix} =: \mathbf{D} \quad .$$

Thus, again $\mathbf{M}^\tau = \mathbf{B}\mathbf{D}^\tau\mathbf{B}^{-1}$, and the claim is that for our example \mathbf{D}^τ reads as

$$\mathbf{D} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & \lambda_2^\tau & \binom{\tau}{1}\lambda_2^{\tau-1} & \binom{\tau}{2}\lambda_2^{\tau-2} & 0 \\ 0 & 0 & \lambda_2^\tau & \binom{\tau}{1}\lambda_2^{\tau-1} & 0 \\ 0 & 0 & 0 & \lambda_2^\tau & 0 \\ 0 & 0 & 0 & 0 & \lambda_3^\tau \end{bmatrix} \quad .$$

All the entries except the ones in the blocks with the binomial coefficient terms are trivial. So let us quickly verify those. For $\tau := 1$ it obviously holds when setting $\binom{\tau}{n} := 0$ for $n > \tau$. So assume the claim holds for $\tau := k$. Then we have

$$\begin{aligned} \mathbf{D}^{k+1} &= \mathbf{D}^k \mathbf{D} \\ &= \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & \lambda_2^k & \binom{k}{1}\lambda_2^{k-1} & \binom{k}{2}\lambda_2^{k-2} & 0 \\ 0 & 0 & \lambda_2^k & \binom{k}{1}\lambda_2^{k-1} & 0 \\ 0 & 0 & 0 & \lambda_2^k & 0 \\ 0 & 0 & 0 & 0 & \lambda_3^k \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & 0 & 0 & 0 \\ 0 & \lambda_2 & 1 & 0 & 0 \\ 0 & 0 & \lambda_2 & 1 & 0 \\ 0 & 0 & 0 & \lambda_2 & 0 \\ 0 & 0 & 0 & 0 & \lambda_3 \end{bmatrix} \quad . \end{aligned}$$

Thus, at $j \geq i$, $d := j - i$, $(\mathbf{D}^k)_{i,j} = \binom{k}{d} \lambda_2^{k-d}$, and hence we have

$$\begin{aligned}
(\mathbf{D}^{k+1})_{i,j} &= (\mathbf{D}^k)_{i,j-1} (\mathbf{D})_{j-1,j} + (\mathbf{D}^k)_{i,j} (\mathbf{D})_{j,j} \\
&= (\mathbf{D}^k)_{i,i+(d-1)} + (\mathbf{D}^k)_{i,i+d} \lambda_2 \\
&= \binom{k}{d-1} \lambda_2^{k-d+1} + \binom{k}{d} \lambda_2^{k-d} \lambda_2 \\
&= \binom{k}{d-1} \lambda_2^{k-d+1} + \binom{k}{d} \lambda_2^{k-d+1} \\
&= \left(\binom{k}{d-1} + \binom{k}{d} \right) \lambda_2^{k-d+1} \\
&\stackrel{\vee}{=} \binom{k+1}{d} \lambda_2^{k+1-d},
\end{aligned}$$

just as expected.

This was just an example, but it is easy to see that we can generalize this, and we get that the absolute value of every entry in \mathbf{D}^τ , except the top left 1, is $\mathcal{O}(|\lambda_2^+|^\tau)$, where λ_2^+ is defined s.t. $|\lambda_2^+| = |\lambda_2| + \epsilon$ for some $\epsilon \in \mathbb{R}_{>0}$. Note that $|\binom{\tau}{d} \lambda_2^{\tau-d}| \in \mathcal{O}(|\lambda_2^+|^\tau)$ for each $d \in \mathbb{N}$.

The rest is trivial, as all the the entries in \mathbf{B} and \mathbf{B}^{-1} are really just constants, and hence when calculating $(\mathbf{M}^\tau)_{ba} = (\mathbf{B} \mathbf{D}^\tau \mathbf{B}^{-1})_{ba}$, we have

$$\begin{aligned}
(\mathbf{B} \mathbf{D}^\tau \mathbf{B}^{-1})_{ba} &= \left(\mathbf{B} \left(\begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix} + \begin{bmatrix} c_{11} & c_{12} & \cdots & c_{1n} \\ c_{21} & c_{22} & \cdots & c_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n1} & c_{n2} & \cdots & c_{nn} \end{bmatrix} \mathcal{O}(|\lambda_2^+|^\tau) \right) \mathbf{B}^{-1} \right)_{ba} \\
&= \boldsymbol{\mu}_b \pm \mathcal{O}(|\lambda_2^+|^\tau),
\end{aligned}$$

for some $c_{ij} \in \mathbb{C}$, $|c_{ij}| = 1$. The rest of the proof is identical to the one given.

1.2 No Markov Chain with Power-Law Behavior

We are interested in cases where $I(X; Y)$ decays to 0, as this is implied by power-law behavior (see definition ?? or definition ??). This means that for increasing τ , we get in the limit $\tau \rightarrow \infty$ that $P(a, b) = P(a)P(b | a) = P(a)(\mathbf{M}^\tau)_{ba} \stackrel{!}{=} P(a)P(b)$, and hence $(\mathbf{M}^\tau)_{ba} = P(b)$ for every $a \in S$. Clearly, this means that \mathbf{M}^τ must converge to a stationary matrix \mathbf{M}_μ , where all columns are equal, at least given the case that $P(b)$ converges for all $b \in S$. But thanks to the following lemma, we can assume that this is the case:

Lemma 1.1 (Convergence of $P(b)$ is Necessary for Power-Law Behavior). *In order for $I(X; Y)$ to converge to 0 for increasing τ , we must have that $\lim_{t \rightarrow \infty} \mathbf{M}^t \mathbf{p}_0 = \boldsymbol{\mu}$ for every $\mathbf{p}_0 \in \Delta$ (and hence $P(b)$ converges).*

Proof. In order for $I(X; Y)$ to converge to 0, we still must have that $|(\mathbf{M}^\tau)_{ba} - P(X_{t_0+\tau} = b)| \xrightarrow{\tau \rightarrow \infty} 0$ for every $a \in S$. But note that for this to happen we must have close to equal columns for increasing τ , because $(\mathbf{M}^\tau)_{ba}$ must be independent of a . But such a matrix must be stationary (and hence $P(b)$ converges). \square

Theorem 1.2 (No Markov Chain with Power-Law Behavior). *There is no Markov chain with weak power-law behavior (and hence also strong power-law behavior).*

Proof. For the sake of contradiction, assume that \mathbf{M} has weak lower bound power-law behavior and decaying behavior (\Leftarrow power-law behavior). Due to the previous results, we know that \mathbf{M}^τ converges to a stationary matrix \mathbf{M}_μ . (The following discussion is facultative; it aims to enhance the reader's understanding). \mathbf{M}_μ must contain 0-entries, as if it didn't, it would mean that \mathbf{M} is irreducible and aperiodic based on corollary ??, and hence \mathbf{M} would have exponential decay.

Thus, let's focus on all the rows of \mathbf{M}_μ with non-0-entries with associated states S_C . We see that for all $i \in S_C, j \in S$ we have that $(\mathbf{M}_\mu)_{ij} > 0$ (especially for $j \in S_C$). Hence, the set S_C describes a closed communication class. Furthermore, there cannot be another closed communication class, since in a closed communication class the transition probabilities between any two states cannot approach 0. Thus, all other states $S \setminus S_C =: S_O$ are in an open communication class.

For indexing reasons, we assume without loss of generality that $S_C = \{1, \dots, n'\} \subseteq S$. Since S_C is a closed communication class, we have that for all $j \in S_C, i \in S_O, t \in \mathbb{N}_{>0} : (\mathbf{M}^t)_{ij} = 0$ (especially for $t := 1$). Thus for all $i, j \in S_C, t \in \mathbb{N}_{>0}$:

$$\begin{aligned} (\mathbf{M}^{t+1})_{ij} &= (\mathbf{M}^t \mathbf{M})_{ij} \\ &= \sum_{k \in S} (\mathbf{M}^t)_{ik} (\mathbf{M})_{kj} \\ &= \sum_{k \in S_C} (\mathbf{M}^t)_{ik} (\mathbf{M})_{kj} \quad . \end{aligned}$$

Now, let \mathbf{M}_{S_C} be the submatrix of \mathbf{M} containing and only containing the transition probability entries for the states in S_C (with our assumption \mathbf{M}_{S_C} is the top left submatrix of \mathbf{M}). From our result, we see that $(\mathbf{M}^t)_{S_C} = (\mathbf{M}_{S_C})^t$.

But $(\mathbf{M}^t)_{S_C}$ converges to a positive matrix, and hence so must $(\mathbf{M}_{S_C})^t$. Hence, \mathbf{M}_{S_C} must be irreducible aperiodic based on corollary ?? again. Thus, $(\mathbf{M}_{S_C})^t$ converges with exponential decay (with a basis of $|\lambda_2| < 1$ of \mathbf{M}_{S_C}), and hence so does $(\mathbf{M}^t)_{S_C}$.

Quickly note that λ_2 of \mathbf{M}_{S_C} is an eigenvalue of \mathbf{M} as well: There must be (at least) one associated eigenvector $\mathbf{v} \in \mathbb{C}^{|S_C|}$ s.t. $\mathbf{M}_{S_C} \mathbf{v} = \lambda_2 \mathbf{v}$. Now, extend \mathbf{v} to

an eigenvector \mathbf{v}' of \mathbf{M} by adding zeros in all the places associated with states in S_O (here: at the end). For $i \in S_C$ we have:

$$\begin{aligned}
(\mathbf{M}\mathbf{v}')_i &= \sum_{j \in S} \mathbf{M}_{ij} \mathbf{v}'_j \\
&\stackrel{\text{0-entries of } \mathbf{v}'}{=} \sum_{j \in S_C} \mathbf{M}_{ij} \mathbf{v}'_j \\
&= \sum_{j \in S_C} \mathbf{M}_{ij} \mathbf{v}_j \\
&\stackrel{i \in S_C}{=} \sum_{j \in S_C} (\mathbf{M}_{S_C})_{ij} \mathbf{v}_j \\
&= (\mathbf{M}_{S_C} \mathbf{v})_i \\
&= \lambda_2 \mathbf{v}_i \stackrel{\check{}}{=} \lambda_2 \mathbf{v}'_i \quad ,
\end{aligned}$$

and for $i \in S_O$ it follows that

$$\begin{aligned}
(\mathbf{M}\mathbf{v}')_i &= \sum_{j \in S} \mathbf{M}_{ij} \mathbf{v}'_j \\
&\stackrel{\text{0-entries of } \mathbf{v}'}{=} \sum_{j \in S_C} \mathbf{M}_{ij} \mathbf{v}'_j \\
&\stackrel{i \in S_O}{=} \sum_{j \in S_C} 0 \cdot \mathbf{v}'_j \\
&= 0 \stackrel{\check{}}{=} \lambda_2 \mathbf{v}'_i \quad .
\end{aligned}$$

The only things left to do are, first, to verify that $|\lambda_2|$ of \mathbf{M}_{S_C} is less than or equal to the absolute value of the second largest eigenvalue of \mathbf{M} ; and second, to show the convergence of \mathbf{M}_{S_O} and that it also exhibits exponential decay (with a base less than or equal to $|\lambda_2|$ of \mathbf{M}). Luckily, we achieve all our goals thanks to the following observation:

(Start of the actual proof.) Since \mathbf{M}^t converges to \mathbf{M}_μ , there must exist an $m \in \mathbb{N}_{>0}$ s.t. \mathbf{M}^m and \mathbf{M}^{m+1} both have a positive row. Thus, thanks to corollary ?? we can apply the Perron-Frobenius Theorem (with the exception that the eigenvector is non-negative instead of positive) to \mathbf{M}^m and \mathbf{M}^{m+1} , and thanks to lemma ?? also to \mathbf{M} itself. Of course, $\lambda_{max} = 1$ with the associated eigenvector μ . (So we know that $|\lambda_2| < 1$ of \mathbf{M}_{S_C} is less than or equal to the absolute value of the second largest eigenvalue of \mathbf{M} .)

Furthermore, we of course have that $\lim_{\tau \rightarrow \infty} (\mathbf{M}^\tau)_{ba} = \mu_b$ ($:= \lim_{\tau \rightarrow \infty} P(b)$). But those were the only two assumptions made in section 1.1.1. Thus, by the same logic, we get $I(X_{t_0}; X_{t_0+\tau}) \in \mathcal{O}(|\lambda_2^+|^\tau)$. \square