

1 Mutual Information in Markov Chains

1.1 Exponential Decay in Irreducible Aperiodic 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

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

where \mathbf{M}_μ is the matrix whose columns all consist of the unique stationary probability distribution μ .

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 $\mathbf{M}^{t_0} \approx \mathbf{M}_\mu$. We will use this "equality" later.

Our goal now is to quantify the mutual information of 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, that is $P'(X, Y) := P(X) \cdot P(Y)$. We use the Kullback-Leibler divergence, so our target expression becomes

$$D(P(X, Y) \parallel P'(X, Y)) \quad .$$

Note that of course this divergence $I(X, Y) := D(P(X, Y) \parallel P'(X, Y))$ 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

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

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.

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:

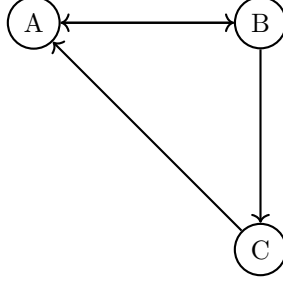


Figure 1: A simple irreducible aperiodic Markov chain. Note that if $X_{t_0} = C$, then we know that $X_{t_0+1} = A$.

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

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 "wash out" over time and converge to the stationary one, 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 "washing out". 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 "washing out" 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 .

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

$$D(P(X, Y) \| P'(X, Y)) = \sum_{(a, b) \in S^2} P(a, b) \log_B \frac{P(a, b)}{P(a)P(b)} \quad .$$

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

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

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 $P(a) \approx \mu_a$ and $P(b) \approx \mu_b$, and also $P(b|a) = (\mathbf{M}^\tau)_{ba}$. Thus:

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

Let us now focus on $(\mathbf{M}^\tau)_{ba}$. For simplicity, we consider the case that \mathbf{M} is diagonalizable. Note that since \mathbf{M} is irreducible and aperiodic, we have that $1 = \lambda_1 > |\lambda_2| \geq \dots \geq |\lambda_n|$. The authors provide proof for the other case as well. But for now, 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 $\tau \rightarrow \infty$? $\boldsymbol{\mu}_b$ of course. What does \mathbf{B} look like? Well, this is very hard to tell, it at least should have a scaled version of $\boldsymbol{\mu}$ in its first column. But we cannot really infer any information about \mathbf{B}^{-1} . But we know

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

So we know that

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

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

This is looking promising, as this means that the discrepancy between $(\mathbf{M}^\tau)_{ba}$ and $\boldsymbol{\mu}_b$ 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. Note that this step deviates from the procedure in the paper (own interpretation, informal!). Thus:

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

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

Remark 1.1. The above proof should also work without the approximation $P(a) \approx \mu_a$, so t_0 musn't be big.

Remark 1.2. Based on the proof, we see that if the distance between M^τ and 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 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 M would be

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

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

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

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

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

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

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

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 $\mathcal{O}(|\lambda_2|^\tau)$. 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

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

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

1.2 Other Cases

We are interested in cases where $I(X, Y)$ decays to 0. This means that for increasing τ , we get in the limit $\tau \rightarrow \infty$ that $P(a, b) = P(a)P(a | b) = P(a)(\mathbf{M}^\tau)_{ba} \stackrel{!}{=} P(a)P(b)$, 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 $b \in S$.

1.2.1 Case: $P(b)$ Does Not Converge

The problem is that the convergence of $P(b)$ might depend on the initial probability distribution vector \mathbf{p}_0 . But actually, note that we still 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), so we cannot have $I(X, Y) = 0$ if $P(b)$ does not converge.

1.2.2 Case: $P(b)$ Does Converge

So now we assume that \mathbf{M}^τ converges to a stationary matrix \mathbf{M}_μ . \mathbf{M}_μ must contain 0-entries, as if it didn't, it would mean that \mathbf{M} is irreducible and aperiodic based on remark ??, and we already have covered that.

So this must mean that the set S_C of all the states defined by the rows of \mathbf{M}_μ that are not 0 form the only closed communication class (and hence all other states are in an open communication class. All those states form the set S_O). Hence, after $n = |S|$ ticks we already have that all the rows associated with states in S_O are filled up with 0-entries.

So we look at \mathbf{M}^n . Imagine for simplicity that the top rows are zeroed out. But then we can *solely* focus on the sub-matrix \mathbf{M}_{S_C} defined by the set of states S_C . But \mathbf{M}_{S_C} converges to a positiv matrix, so \mathbf{M}_{S_C} must be aperiodic based on remark ?? again. This means that the irreducible aperiodic matrix \mathbf{M}_{S_C} converges with exponential decay, and hence so does \mathbf{M}^n . Furthermore, we also know that \mathbf{M} itself converges, and hence it also must experience exponential decay towards \mathbf{M}_μ .