

Maximum Kolmogorov-Sinai Entropy Versus Minimum Mixing Time in Markov Chains

M. Mihelich¹ · B. Dubrulle¹ · D. Paillard² ·
Q. Kral^{3,4} · D. Faranda²

Received: 4 December 2015 / Accepted: 6 September 2017 / Published online: 21 November 2017
© Springer Science+Business Media, LLC 2017

Abstract We establish a link between the maximization of Kolmogorov Sinai entropy (KSE) and the minimization of the mixing time for general Markov chains. Since the maximisation of KSE is analytical and easier to compute in general than mixing time, this link provides a new faster method to approximate the minimum mixing time dynamics. It could be interesting in computer sciences and statistical physics, for computations that use random walks on graphs that can be represented as Markov chains.

Keywords Entropy · Markov chain · Mixing time

Many modern techniques of physics, such as computation of path integrals, now rely on random walks on graphs that can be represented as Markov chains. Techniques to estimate the number of steps in the chain to reach the stationary distribution (the so-called “mixing time”), are of great importance in obtaining estimates of running times of such sampling algorithms [1] (for a review of existing techniques, see e.g. [2]). On the other hand, studies of the link between the topology of the graph and the diffusion properties of the random walk on this graph are often based on the entropy rate, computed using the Kolmogorov-Sinai entropy (KSE) [3,4]. For example, one can investigate dynamics on a network maximizing the KSE to study optimal diffusion [3], or obtain an algorithm to produce equiprobable paths on non-regular graphs [5].

✉ M. Mihelich
ma.mihelich@gmail.com

¹ Laboratoire SPHYNX, Service de Physique de l’Etat Condensé, DSM, CEA Saclay, CNRS UMR 3680, 91191 Gif-sur-Yvette, France

² Laboratoire des Sciences du Climat et de l’Environnement, IPSL, CEA-CNRS-UVSQ, UMR 8212, 91191 Gif-sur-Yvette, France

³ Institute of Astronomy, University of Cambridge, Madingley Road, Cambridge CB3 0HA, UK

⁴ LESIA-Observatoire de Paris, UPMC University Paris 06, University Paris-Diderot, 92195 Meudon Cedex, France

In this letter, we establish a link between these two notions by showing that for a system that can be represented by Markov chains, *a non trivial relation exists between the maximization of KSE and the minimization of the mixing time*. Since KSE are easier to compute in general than mixing time, this link provides a new faster method to approximate the minimum mixing time that could be interesting in computer sciences and statistical physics and gives a physical meaning to the KSE. We first show that on average, the greater the KSE, the smaller the mixing time, and we correlated this result to its link with the transition matrix eigenvalues. Then, we show that the dynamics that maximises KSE is close to the one minimizing the mixing time, both in the sense of the optimal diffusion coefficient and the transition matrix.

Consider a network with m nodes, on which a particle jumps randomly. This process can be described by a finite Markov chain defined by its adjacency matrix A and its transition matrix P . $A(i, j) = 1$ if and only if there is a link between the nodes i and j and 0 otherwise. $P = (p_{ij})$ where p_{ij} is the probability for a particle in i to hop on the j node. Let us introduce the probability density at time n $\mu_n = (\mu_n^i)_{i=1\dots m}$ where μ_n^i is the probability that a particle is at node i at time n . Starting with a probability density μ_0 , the evolution of the probability density writes: $\mu_{n+1} = P^t \mu_n$ where P^t is the transpose matrix of P .

Within this paper, we assume that the Markov chain is irreducible and thus has a unique stationary state.

Let us define:

$$d(n) = \max \| (P^t)^n \mu - \mu_{stat} \| \forall \mu, \tag{1}$$

where $\| \cdot \|$ is a norm on \mathbb{R}^n . For $\epsilon > 0$, the mixing time, which corresponds to the time such that the system is within a distance ϵ from its stationary state is defined as follows:

$$t(\epsilon) = \min n, d(n) \leq \epsilon. \tag{2}$$

For a Markov chain the KSE takes the analytical form [6]:

$$h_{KS} = - \sum_{ij} \mu_{stat_i} p_{ij} \log(p_{ij}). \tag{3}$$

Random m size Markov matrices are generated by assigning to each p_{ij} a random number between 0 and 1 and by normalized each row. The mean KSE is plotted versus the mixing time (Fig. 1) by working out h_{KS} and $t(\epsilon)$ for each random matrix. (Fig. 1) shows that KSE is on average a decreasing function of the mixing time.

We stress the fact that this relation is only true on average. We can indeed find two special Markov chains $P1$ and $P2$ such that $h_{KS}(P1) \leq h_{KS}(P2)$ and $t_1(\epsilon) \leq t_2(\epsilon)$. We illustrate this point further.

The link between the mixing time and the KSE can be understood via their dependence as a function of the transition matrix eigenvalues. More precisely, we have found a heuristic connection between the second largest eigenvalue of the transition matrix and the KSE. A general irreducible transition matrix P is not necessarily diagonalizable on \mathbb{R} . However, since P is chosen randomly, it is almost everywhere diagonalizable on \mathbb{C} . According to Perron Frobenius theorem, the largest eigenvalue is 1 and the associated eigen-space is one-dimensional and equal to the vectorial space generated by μ_{stat} . Without loss of generality, we can label the eigenvalues in decreasing order of their module:

$$1 = \lambda_1 > |\lambda_2| \geq \dots \geq |\lambda_m| \geq 0$$

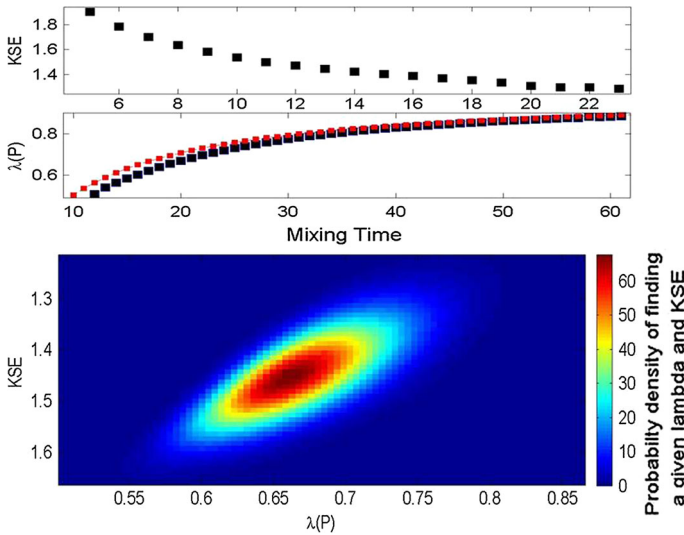


Fig. 1 Averaged KSE versus mixing time for 10^6 random $m = 10$ size matrices and averaged $\lambda(P)$ versus mixing time for 10^6 random $m = 10$ size matrices in curve blue and $f(t) = \epsilon^{1/t}$ in red (top). Heat map representing KSE and $\lambda(P)$ (bottom) with in z axis the probability density of finding a given KSE and a given $\lambda(P)$. Here $\epsilon = 10^{-3}$ and the norm is chosen to be the euclidian one

The convergence speed toward μ_{stat} is given by the second maximum module of the eigenvalues of P [7, 8]:

$$\lambda(P) = \max_{i=2\dots m} |\lambda_i| = |\lambda_2|$$

Moreover, it is well known that $\lambda(P) \propto \epsilon^{1/t(\epsilon)}$. Hence, the smaller $\lambda(P)$ the shorter the mixing time (Fig. 1). h_{KS} being a decreasing function of $t(\epsilon)$ and $\lambda(P)$ being an increasing function of $t(\epsilon)$, we deduce that h_{KS} is a decreasing function of $\lambda(P)$.

This result can be demonstrated anatically for 2×2 matrices which have the form $P = \begin{pmatrix} 1-a & a \\ b & 1-b \end{pmatrix}$ and where a and b are taken uniformly in $[0, 1]^2$ such as $a + b \leq 1$. In fact, the average KS entropy function of $\lambda(P)$ takes the form:

$$\begin{aligned} E_{KS}(\lambda(P)) &= \frac{1}{3}(1 - \lambda(P)) \log(1 - \lambda(P)) - \frac{5}{18} * (1 - \lambda(P)) \\ &+ \frac{1}{3} * \frac{\lambda(P)^3 \log(\lambda(P))}{(1 - \lambda(P))^2} + \frac{5}{36} \frac{(\lambda(P)^2 + \lambda(P) - 4/5)}{1 - \lambda(P)} \end{aligned} \tag{4}$$

With little algebra one can show that the derivative is strictly negative on $]0, 1[$ thus the function is strictly decreasing. This result can be numerically verified by drawing random matrices et by calculating $\lambda(P)$ and E_{KS} .

This link between maximum KSE and minimum mixing time actually also extends naturally to optimal diffusion coefficients. Such a notion has been introduced by Gomez-Gardenes and Latora [3] in networks represented by a Markov chain depending on a diffusion coefficient. Based on the observation that in such networks, KSE has a maximum as a function of the diffusion coefficient, they define an optimal diffusion coefficient as the value of the diffusion corresponding to this maximum. In the same spirit, one could compute an opti-

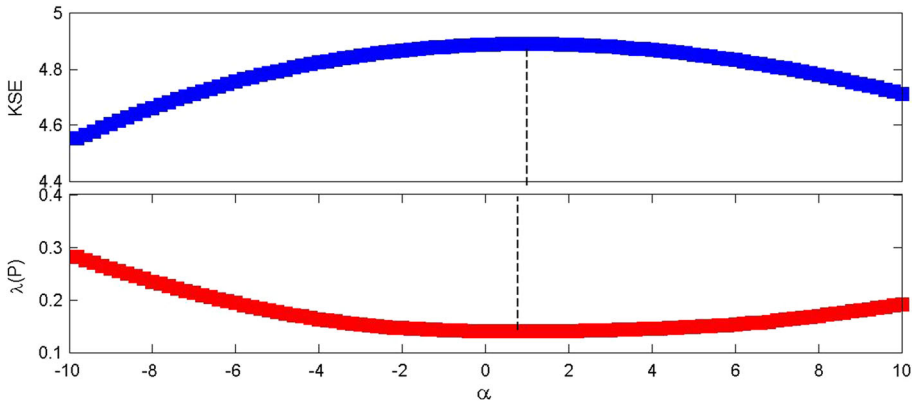


Fig. 2 KSE (top) and $\lambda(P)$ (bottom) function of α for a network of size $m = 400$ with a proportion of 0 in A equal to $1/3$

mal diffusion coefficient with respect to the mixing time, corresponding to the value of the diffusion coefficient which minimizes the mixing time -or equivalently the smallest second largest eigenvalue $\lambda(P)$. This would roughly correspond to the diffusion model reaching the stationary time in the fastest time. To define such an optimal diffusion coefficient, we follow Gomez and Latora and vary the transition probability depending on the degree of the graph nodes. More accurately, if $k_i = \sum_j A(i, j)$ denotes the degree of node i , we set:

$$p_{ij} = \frac{A_{ij}k_j^\alpha}{\sum_j A_{ij}k_j^\alpha}. \tag{5}$$

If $\alpha < 0$ we favor transitions towards low degrees nodes, if $\alpha = 0$ we find the typical random walk on network and if $\alpha > 0$ we favor transitions towards high degrees nodes. We assume here that A is symmetric. It may then be checked that the stationary probability density is equal to:

$$\pi_{stat_i} = \frac{c_i k_i^\alpha}{\sum_j c_j k_j^\alpha}, \tag{6}$$

where $c_i = \sum_j A_{ij}k_j^\alpha$,

Using Eqs. (5) and (6), we check that the transition matrix is reversible and then has m real eigenvalues. From this stationary probability density, we can thus compute both the KSE and the second largest eigenvalue $\lambda(P)$ as a function of α . The result is provided in (Fig. 2).

We observe in (Fig. 2) that the KS entropy has a maximum at a value that we denote α_{KS} , in agreement with the findings of [3]. Likewise, $\lambda(P)$ (i.e. the mixing time) presents a minimum for $\alpha = \alpha_{mix}$. Moreover, α_{KS} and α_{mix} are close. This means that the two optimal diffusion coefficients are close to each other. Furthermore, looking at the ends of the two curves, we can find two special Markov chains $P1$ and $P2$ such that $h_{KS}(P1) \leq h_{KS}(P2)$ and $t_1(\epsilon) \leq t_2(\epsilon)$, illustrating that the link between KSE and the minimum mixing time is only true in a general statistical sense.

We have thus shown that, for a given transition matrix P (or equivalently for given jump rules) the greater the KSE, the smaller the mixing time. We now investigate whether a similar property holds for dynamics, i.e. whether transition rules that maximise KSE are close to the ones minimizing the mixing time. For a given network, i.e. for a fixed A , there is a well known procedure to compute the transition matrix P_{KS} which maximizes the KSE with the

constraints $A(i, j) = 0 \Rightarrow P_{KS}(i, j) = 0$ [5]. It proceeds as follow: let us note λ the greatest eigenvalue of A and Ψ the normalized eigenvector associated i.e $A\Psi = \lambda\Psi$ and $\sum_i \Psi_i^2 = 1$. We define P_{KS} such that:

$$P_{KS}(i, j) = \frac{A(i, j)}{\lambda} \frac{\Psi_j}{\Psi_i}. \tag{7}$$

We have $\forall i \sum_j P_{KS}(i, j) = 1$. Moreover, using the fact that A is symmetric we find:

$$\sum_j P_{KS}(j, i) \Psi_j^2 = \sum_j \frac{A(j, i) \Psi_i \Psi_j}{\lambda} = \Psi_i^2. \tag{8}$$

Hence, $P_{KS}^t \Psi^2 = \Psi^2$ and the stationary density of P_{KS} is $\pi_{stat} = \Psi^2$.

Using Eqs. (3) and (7), we have:

$$h_{KS} = -\frac{1}{\lambda} \sum_{(i,j)} A(i, j) \Psi_i \Psi_j \log \left(\frac{A(i, j)}{\lambda} \frac{\Psi_j}{\Psi_i} \right). \tag{9}$$

Eq. (9) can be split in two terms:

$$\begin{aligned} h_{KS} &= \frac{1}{\lambda} \sum_{(i,j)} A(i, j) \Psi_i \Psi_j \log(\lambda) \\ &\quad - \frac{1}{\lambda} \sum_{(i,j)} A(i, j) \Psi_i \Psi_j \log \left(A(i, j) \frac{\Psi_j}{\Psi_i} \right). \end{aligned} \tag{10}$$

The first term is equal to $\log(\lambda)$ because Ψ is an eigenvector of A and the second term is equal to 0 due to the symmetry of A . Thus:

$$h_{KS} = \log(\lambda). \tag{11}$$

Moreover, for a Markov chain the number of trajectories of length n is equal to $N_n = \sum_{(i,j)} (A^n)(i, j)$. For a Markov chain the KSE can be seen as the time derivative of the path entropy leading that KSE is maximal when the paths are equiprobable. For an asymptotic long time the maximal KSE is:

$$h_{KS_{max}} = \frac{\log(N_n)}{n} \rightarrow \log(\lambda), \tag{12}$$

by diagonalizing A . Using Eqs. (11) and (12) we find that P_{KS} defined as in Eq. (7) maximises the KSE. Finally P_{KS} verifies $\pi_{stat_i} P_{KS}(i, j) = \pi_{stat_j} P_{KS}(j, i) \forall (i, j)$ and thus P_{KS} is reversible.

In a similar way, we can search for a transition matrix P_{mix} which minimizes the mixing time -or, equivalently the transition matrix minimizing its second eigenvalue $\lambda(P)$. This problem is much more difficult to solve than the first one, given that the eigenvalues of P_{mix} can be complex. Nevertheless, two cases where the matrix P_{mix} is diagonalizable on \mathbb{R} can be solved [7]: the case where P_{mix} is symmetric and the case where P_{mix} is reversible for a given fixed stationary distribution. Let us first consider the case where P is symmetric. The minimisation problem takes the following form:

$$\left\{ \begin{array}{l} \min_{P \in S_n} \lambda(P) \\ P(i, j) \geq 0, P * \mathbf{1} = \mathbf{1} \\ A(i, j) = 0 \Rightarrow P(i, j) = 0 \end{array} \right. \tag{13}$$

given the strict convexity of λ and the compactness of the stochastic matrices, this problem admits an unique solution.

P is symmetric thus $\mathbf{1}$ is an eigenvector associated with the largest eigenvalue of P . Then the eigenvectors associated to $\lambda(P)$ are in the orthogonal of $\mathbf{1}$. The orthogonal projection on $\mathbf{1}^\perp$ writes: $I_d - \frac{1}{n}\mathbf{1}\mathbf{1}^t$

Moreover, if we take the matrix norm associated with the euclidian norm i.e. for M any matrix $|||M||| = \max \frac{||MX||_2}{||X||_2} X \in \mathbb{R}^n X \neq 0$ it is equal to the square root of the largest eigenvalue of MM^t and then if M is symmetric it is equal to $\lambda(M)$.

Then the minimization problem can be rewritten:

$$\begin{cases} \min_{P \in S_n} |||(I_d - \frac{1}{n}\mathbf{1}\mathbf{1}^t)P(I_d - \frac{1}{n}\mathbf{1}\mathbf{1}^t)||| = |||P - \frac{1}{n}\mathbf{1}\mathbf{1}^t||| \\ P(i, j) \geq 0, P * \mathbf{1} = \mathbf{1} \\ A(i, j) = 0 \Rightarrow P(i, j) = 0 \end{cases} \tag{14}$$

We solve this constrained optimization problem (Karush-Kuhn-Tucker) with Matlab and we denote P_{mix} the matrix which minimizes this system.

We remark that the mixing time of P_{KS} is smaller than the mixing time of P_{mix} . This is coherent because in order to calculate P_{KS} we take the minimum on all the matrix space whereas to calculate P_{mix} we restrict us to the symmetric matrix space. Nevertheless, we can go a step further and calculate, the stationary distribution being fixed, the reversible matrix which minimizes the mixing time. If we note π the stationary measure and $\Pi = diag(\pi)$. Then P is reversible if and only if $\Pi P = P^t \Pi$. Then in particular $\Pi^{\frac{1}{2}} P \Pi^{-\frac{1}{2}}$ is symmetric and has the same eigenvalues as Π . Finally, $p = (\sqrt{\pi_1}, \dots, \sqrt{\pi_n})$ is an eigenvector of $\Pi^{\frac{1}{2}} P \Pi^{-\frac{1}{2}}$ associated to the eigenvalue 1. Then the minimization problem can be written as the following system:

$$\begin{cases} \min_P |||(I_d - \frac{1}{n}\mathbf{q}\mathbf{q}^t)\Pi^{\frac{1}{2}} P \Pi^{-\frac{1}{2}}(I_d - \frac{1}{n}\mathbf{q}\mathbf{q}^t)||| \\ = |||\Pi^{\frac{1}{2}} P \Pi^{-\frac{1}{2}} - \frac{1}{n}\mathbf{q}\mathbf{q}^t||| \\ P(i, j) \geq 0, P * \mathbf{1} = \mathbf{1}, \Pi P = P^t \Pi \\ A(i, j) = 0 \Rightarrow P(i, j) = 0 \end{cases} \tag{15}$$

When we implement this problem in Matlab with $\pi = \pi_{KS}$ we find a matrix P_{mix} such that naturally $\lambda(P_{mix}) \leq \lambda(P_{KS})$. Moreover we can compare both dynamics by evaluating $|||P_{KS} - P_{mix}|||$ compared to $|||P_{KS}|||$ which is approximatively equal to $|||P_{mix}|||$. We remark that $|||P_{KS} - P_{mix}|||$ depends on the density ρ of 0 in the matrix A . For a density equal to 0 the matrices P_{KS} and P_{mix} are equal and the quantity $|||P_{KS} - P_{mix}|||$ will increase continuously when ρ increases. This is shown in (Fig. 3).

From this, we conclude that the rules which maximize the KSE are close to those which minimize the mixing time. This becomes increasingly accurate as the fraction of removed links in A is weaker. Since the calculation of P_{mix} quickly becomes tedious for quite large values of m , we offer here a much cheaper alternative by computing P_{KS} instead of P_{mix} .

In the previous study we fixed the dynamics i.e the adjacency matrix A . We can ask ourself what happens when we impose a prescribed stationary distribution on the graph and maximize its KS entropy. If we impose a stationnary measure $\mu_{stat} = (\mu_i)$ and we leave the dynamics totally free, we can prove analytically from Dixit work [9, 10] that the matrice which maximizes KS entropy is exactly the same which minimizes the mixing time and where the coefficients of the matrix are $p_{ij} = \mu_j$.

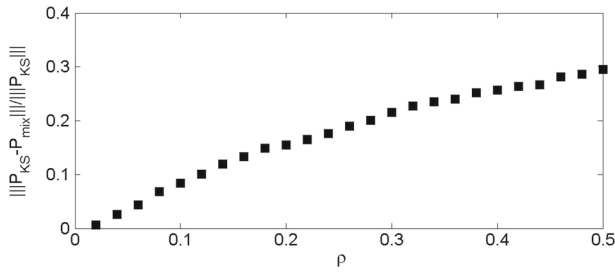


Fig. 3 $\|P_{KS} - P_{mix}\| / \|P_{KS}\|$ as a function of the density ρ of 0 present in A

Moreover, maximizing the KSE appears today as a method to describe out of equilibrium complex systems [11], to find natural behaviors [5] or to define optimal diffusion coefficients in diffusion networks. This general observation however provides a possible rationale for selection of stationary states in out-of-equilibrium physics: it seems reasonable that in a physical system with two simultaneous equiprobable possible dynamics, the final stationary state will be closer to the stationary state corresponding to the fastest dynamics (smallest mixing time). Through the link found in this letter, this state will correspond to a state of maximal KSE. If this is true, this would provide a more satisfying rule for selecting stationary states in complex systems such as climate than the maximization of the entropy production, as already suggested in [12].

Acknowledgements Martin Mihelich thanks IDEEX Paris-Saclay for financial support. Quentin Kral was supported by the French National Research Agency (ANR) through contract ANR-2010 BLAN-0505-01 (EXOZODI).

References

1. Bhakta, P., Miracle, S., Randall, D., Streib, A.P.: Mixing times of markov chains for self-organizing lists and biased permutations. In: Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, pp. 1–15. SIAM (2013)
2. Guruswami, V.: Rapidly mixing markov chains: a comparison of techniques. <http://cs.washington.edu/homes/venkat/pubs/papers.html> (2000)
3. Gómez-Gardeñes, J., Latora, V.: Entropy rate of diffusion processes on complex networks. *Phys. Rev. E* **78**(6), 065102 (2008)
4. Ochab, J.K.: Static and dynamic properties of selected stochastic processes on complex networks. PhD thesis, Institute of Physics (2013)
5. Burda, Z., Duda, J., Luck, J.M., Waclaw, B.: Localization of the maximal entropy random walk. *Phys. Rev. Lett.* **102**(16), 160602 (2009)
6. Billingsley, P.: Ergodic Theory and Information. Wiley, New York (1965)
7. Boyd, Stephen, Diaconis, Persi, Xiao, Lin: Fastest mixing markov chain on a graph. *SIAM Rev.* **46**(4), 667–689 (2004)
8. Bremaud, P.: Markov Chains: Gibbs fields, Monte Carlo Simulation, and Queues, vol. 31. Springer, New York (1999)
9. Dixit, P.D., Dill, K.A.: Inferring microscopic kinetic rates from stationary state distributions. *J. Chem. Theory Comput.* **10**(8), 3002–3005 (2014)
10. Dixit, P.D., Jain, A., Stock, G., Dill, K.A.: Inferring transition rates of networks from populations in continuous-time markov processes. *J. Chem. TheoryComput.* **11**(11), 5464–5472 (2015)
11. Monthus, C.: Non-equilibrium steady states: maximization of the Shannon entropy associated with the distribution of dynamical trajectories in the presence of constraints. *J. Stat. Mech.* **2011**, P03008 (2011)
12. Mihelich, Martin, Dubrulle, Bérengère, Paillard, Didier, Herbert, Corentin: Maximum entropy production versus kolmogorov-sinai entropy in a constrained asep model. *Entropy* **16**(2), 1037–1046 (2014)