

**COLT 2015 (author)**[Help](#) [Log out](#)

My Submissions	COLT 2015	Help questions	Alerts	EasyChair
----------------	-----------	----------------	--------	-----------

## Author Response Information for Submission 167

Response Letter(s)

Response Letter	
Response:	<p>We thank the referees for the helpful comments. A detailed response follows.</p> <p>Empirical bounds:</p> <p>Empirical bounds are essential for many real learning applications where one must make decisions based on a single sample path, with no other prior information, and hence are of independent interest. Such bounds were also sought by McDonald et al (2011) and Meir (2000) but not achieved. We also think obtaining empirical bounds for <math>\gamma_*</math> is very natural problem.</p> <p>The confidence interval from Thm.4 is fully empirical, and do not require apriori knowledge of <math>\gamma_*</math> or <math>\pi_*</math>. In particular, the interval is always correct even if <math>\gamma_*</math> and <math>\pi_*</math> are actually zero. When <math>\gamma_*</math> and <math>\pi_*</math> are positive, the width of the interval shrinks at a rate shown in Thm.4 (which depends inversely on <math>\gamma_*</math> and <math>\pi_*</math>) and hence we obtain nontrivial confidence intervals.</p> <p>A confidence region for <math>\gamma_*</math> that one can obtain from Thm.1 is a union of two intervals: one around 0, and another around <math>\gamma_*</math>. Call these <math>I_0</math> and <math>I_1</math>, respectively. The width of each shrinks (at a better rate than Thm.4), but <math>\gamma_*=0</math> is never precluded.</p> <p>However, using <math>\text{union}(I_0, I_1)</math> intersected with the Alg.1 interval (a <math>1-2\delta</math> confidence region) eventually precludes <math>\gamma_*=0</math> (and eventually <math>I_0</math>), upon which we are left with a single interval as tight as <math>I_1</math> alone.</p> <p>Nonreversible chains:</p> <p>We discuss an <math>\exp(d)</math> approach to non-reversible chains (Sec.6), leaving <math>\text{poly}(d)</math> as an open problem. The reversible case is important (it arises in MCMC) and also the natural first step in this line of work.</p> <p>E-norms:</p> <p>Indeed, they are mixed-up, and there is a small gap in the submission, which is easily fixed. Upper bound the norm of <math>E_{\{\pi, i\}}</math> by the maximum of the norms of <math>E_{\{\pi, 1\}}</math> and <math>E_{\{\pi, 2\}}</math>. Now, use <math>\max( \sqrt{x}-1 ,  \sqrt{1/x}-1 ) \leq 1/2 \max( x-1 ,  1/x-1 )</math>. The only change in the rest is to redefine <math>\hat{\rho}(\delta)</math> to be <math>2 \max(\rho(\delta), \rho'(\delta)) + 4 \max(\rho(\delta), \rho'(\delta))^2</math>. This does not impact the result's asymptotics, as <math>\rho(\delta)</math> and <math>\rho'(\delta)</math> are of the same order-of-magnitude.</p> <p>Delta:</p>

	<p><math>\tilde{\Delta}</math> is closed-form (easy to analyze); computing <math>\hat{\Delta}</math> requires solving an optimization problem. <math>\tilde{\Delta}</math> is not always smaller. The opposite is true whp: the true <math>P_{\{i,j\}}</math> is a feasible solution, and its deviation from <math>\hat{P}_{\{i,j\}}</math> is bounded by <math>\tilde{\Delta}</math>.</p> <p>Thms.<math>\{1,2,3\}</math>:</p> <p>Bounds are not tight: the polynomial dependence on <math>1/\pi_*</math> and <math>1/\gamma_*</math> do not match (linear in lowerbound, higher degree poly in upperbound). Thm.2 implies lowerbound for both multiplicative and additive accuracy; Thm.3 gives lowerbound for multiplicative accuracy.</p> <p>Bernstein/Samson/Paulin:</p> <p>On page 8, we use standard iid Bernstein (exploiting Markov assumption). However, we also use Theorem 3.8 of Paulin (2015) to prove Theorem 1 (we inadvertently submitted a version that uses looser bounds). Neither Samson/Paulin directly gives the result needed matrices, nor empirical bounds; they are related and will be cited in the revision.</p> <p>More for Rev.3:</p> <p>Standard binomial lowerbounds imply <math>1/\pi_*</math> is necessary to estimate <math>\pi_*</math> multiplicatively accurately.</p> <p>Besides Sym(L) trick, we have another new trick to estimate <math>\pi</math>.</p>
<b>Time:</b>	Apr 14, 17:59 GMT
<b>Letter:</b>	<p>Dear author,</p> <p>Thank you for your submission to COLT 2015. The author response period will be between now and April 15th, Midnight EST.</p> <p>* Note that this is a strict deadline.</p> <p>During this time, you will have access to the current state of your reviews and have the opportunity to submit a response of up to 500 words. Please keep in mind the following during this process:</p> <p>* You can upload your response <b>**only once**</b>, so please edit carefully before you send off.</p> <p>* The response must focus on any factual errors in the reviews and</p>

any questions posed by the reviewers. It must not provide new research results or reformulate the presentation. Try to be as concise and to the point as possible.

\* The author response period is an opportunity to react to the reviews, but not a requirement to do so. Thus, if you feel the reviews are accurate and the reviewers have not asked any questions, then you should not respond.

\* The reviews are as submitted by the PC members, without any coordination between them. Thus, there may be inconsistencies. Furthermore, these are not the final versions of the reviews. The reviews will be updated to take into account discussions by the program committee, and we may find it necessary to solicit other outside reviews after the author response period.

\* The program committee will read your responses carefully and take this information into account during the discussions. On the other hand, the program committee will not directly respond to your responses, neither directly nor in the final versions of the reviews.

The reviews on your paper are attached to this letter. To submit your response you should log on the EasyChair Web site for COLT 2015 and select your submission on the menu.

Best wishes,

Peter Grunwald & Elad Hazan  
COLT 2015 Program Chairs

[\*REVIEWS\*]

<b>Time:</b>	Apr 10, 20:58 GMT
--------------	-------------------

## Reviews

### Review 1

Summary of the paper:

This paper considers the problem of estimating the mixing time in Markov chains when one only has access to a single trajectory. For reversible chains, confidence intervals on the related spectral gap

	<p>are provided that depend on the length of the trajectory, the spectral gap itself, the minimal value of the stationary distribution, as well as the size of the state space. This is complemented by lower bounds that show that dependence on these parameters is necessary.</p> <p>The considered problem is definitely of general interest. However, the main restriction in my view is that the results only hold for reversible Markov chains. The paper mentions several possible applications, but does not discuss in which of them bounds for reversible chains would be sufficient.</p>
Significance and scope:	<p>Moreover, the fact that the size of the confidence intervals in the main theorem depends on the value to estimate as well as some random quantities seriously confines the applicability, as there is no direct access to the size of the confidence interval and therefore to the quality of the estimate. Thus, overall the results seem to be only of theoretical interest.</p>
Novelty and related work:	<p>As far as I can tell, the paper makes a reasonable contribution. However, maybe the result of Theorem 1 (which the authors seem to consider not that important) may be of more interest than Theorem 4, as the confidence intervals for the latter contain some random terms.</p> <p>There are a few issues with soundness and clarity of the paper:</p>
Soundness:	<p>- First of all, it was not clear to me what is the difference between Theorem 1 and 4. For the former, it is claimed that it is not directly usable, since one cannot rule out that the value of <math>\pi^*</math> (the minimal value of the stationary distribution) is 0. However, I don't see why the computation of Algorithm 1 and the respective results of Theorem 4 are better in that respect. In general, the paper assumes that the spectral gap and the minimal stationary probability are positive, but it isn't discussed whether this is a serious restriction.</p> <p>- Similarly, the introduction and comparison of the two different values <math>\hat{\Delta}</math> and <math>\tilde{\Delta}</math> is rather confusing. If the results in Theorem 4 are for the tighter <math>\hat{\Delta}</math>, why is <math>\tilde{\Delta}</math> used in the algorithm? Also, I don't see why <math>\hat{\Delta}</math> is tighter. Since <math>\hat{\Delta}</math> has an additional factor of <math>\sqrt{N_i}</math> in the numerator, at least for big <math>N_i</math> the value of <math>\tilde{\Delta}</math> should be the smaller one. In general, why care about a worse estimate anyway?</p> <p>- The confidence interval of Algorithm 1 / Theorem 4 contains the random values of <math>N_i</math> and <math>\hat{P}_{ij}</math>, which makes the results in my view less applicable than those of Theorem 1. Moreover, the argument concerning the width of the confidence interval is rather confusing. First it is said that one considers the case that <math>n</math> goes to infinity, but then the bounding terms still contain <math>n</math>. Also, for the</p>

$\hat{\Delta}$ -term, I did not see what happened to the  $\sqrt{N_i}$ -term in the numerator when  $n$  goes to infinity. Last but not least, in the final step, it is not clear to me why the bounding term for the product of  $\hat{\rho}(\delta)$  and the  $\sqrt{\cdot}$ -term should be the sum (and not the product) of the individual bounding terms.

- Overall, the proof of Theorem 4 (I didn't check the proof of Theorem 1 in the appendix) could sometimes use a bit more information in some steps, like definition of notation, references to applied results, or some intermediate steps, cf. comments below.

- Finally, there also seems to be a problem with the estimation of the norm of the E-terms on p.10. First, in the bounds for  $\|E_{\pi,1}\|$  and  $\|E_{\pi,2}\|$ , it seems that the two terms are mixed up, that is, for  $\|E_{\pi,1}\|$  one should have the fraction  $\pi/\hat{\pi}$  instead of  $\hat{\pi}/\pi$  and vice versa for  $\|E_{\pi,2}\|$ . More seriously, it seems that the first inequality does not hold in general. If e.g.  $a := \hat{\pi}/\pi = 1/4$ , then  $|\sqrt{a}-1| = 1/2 > 1/2 * |a-1| = 3/8$ .

\*Further comments and corrections\*

- p.2f: In the collection of the main results, some quantities ( $n$ ,  $\delta$ ) are not explained. Also, the algorithm uses some notation like  $A^\#$  that is only introduced later.

- p.6: In Theorem 3, the statement "For any initial state there exists a Markov chain ..." makes no sense, since the states usually do not exist independent of the Markov chain.

- p.6: In Theorem 4, I didn't understand what event the phrase "on the same event" refers to. Further, if  $n$  goes to infinity, it doesn't make sense to still have terms of  $n$  in the formula, so " $n \rightarrow \infty$ " should probably be dropped.

- p.7: In l.4 of the algorithm, 'P' should be ' $\hat{P}$ '.

- p.7: It is claimed that "the interval from Theorem 1 does not include zero", however this contradicts what has been claimed on p.5 in the discussion of Theorem 1.

- p.8: The application of the Bernstein bound should be explained in more detail.

- p.9: For the definition and the results (uniqueness) for the group inverse a reference would be appropriate.

- p.9: In eq.(6), the notation  $[\dots]_+$  is used without explanation.

- p.9: A reference to Weyl's equality wouldn't hurt here.

	<p>- p.10: In the decomposition of <math>L - \hat{L}</math> it should be said that one also uses the definitions of <math>L</math> and <math>\hat{L}</math>. Also in the next step, it seems one uses that <math>  \hat{L}   = 1</math>, which should be mentioned as well.</p> <p>- p.10: In the bound for <math>  E_P  </math>, the exponents for <math>\hat{\pi}_i</math> and <math>\pi_i</math> are missing.</p> <p>While the paper considers an interesting problem, the results only hold in the rather confined setting of reversible Markov chains and seem to be of little practical relevance. Further, there are some problems with soundness and clarity.</p>
Summary of review:	
Significance:	3: (fair)
Novelty and related work:	4: (good)
Soundness:	2: (poor)
Presentation:	3: (fair)
Overall evaluation:	<b>4</b> : (reject: OK paper, but not good enough: issues with technical content, significance or originality.)
Reviewer's confidence:	<b>3</b> : (high)

### Review 2

Summary of the paper:	The work provides estimators for the mixing time of finite and reversible Markov chains from a single sampling sequence. The main result is a fully empirical confidence interval for the spectral gap.
Significance and scope:	The result is clean and solves a clearly defined and important problem about Markov chain estimation: that of providing confidence bound that do not depend on unknown coefficients.
Novelty and related work:	The proofs are clean and use some recent matrix perturbation theory results.
Soundness:	The proofs are detailed and the paper is generally self-consistent.
Summary of review:	A cleanly written paper solving an interesting problem of Markov chain estimation. Some comments about the necessity of the inverse dependence on $\pi^*$ would be nice to have.
Significance:	4: (good)
Novelty and related work:	4: (good)
Soundness:	4: (good)
Presentation:	4: (good)
Overall evaluation:	<b>7</b> : (accept: good paper)
Reviewer's confidence:	<b>1</b> : (low)

### Review 3

## Summary of the paper:

This paper proposes and analyzes a neat way of estimating the mixing time of a reversible (finite state) Markov chain in a fully data-dependent way by means of a single path dynamic of the chain. Upper and lower bounds are provided on the spectral gap and the minimum probability of the stationary distribution, the two quantities which are known to determine mixing rate.

Neatly presented and well-motivated paper on a relevant topic. I really enjoyed the argument that symmetrizes the matrix whose eigenvalues have to be estimated. I expect the argument of using a single run to be of practical relevance.

## Significance and scope:

Comments/requests for clarification:

1) your lower bounds in Sect. 4.2 do not say anything specific about the effort to estimate  $\pi^*$ .

Are the upper bounds in Thm 1 any tight? Why do you need constant  $C$ , i.e., why is error in estimating  $\gamma^*$  additive while the one for  $\pi^*$  multiplicative?

2) Page 8, "by Bernstein inequality..." : which Bernstein ineq. are you using exactly here on  $N_{\{i,j\}}$ ? Pls clarify.

Minor comments (some of which are just typos) :

- page 2: "is assumes"
- page 2, last line: I'd say that  $n$  is the length of the sequence.
- page 4, def of  $\{\hat{\pi}_i\}$ , it should be " $i \in [d]$ "
- page 6, "in he"
- page 7, Item 4 in Algorithm 1:  $P \rightarrow \{\hat{P}\}$

From a technical standpoint, this paper is really hinging on replacing  $P$  with  $\text{sym}(L)$ , then using tighter perturbation schemes, and replacing concentration bounds (obtained by standard blocking) by empirical versions thereof. Still, the result is a neat piece of work.

There are a number of papers on the web (e.g., arxiv) providing concentration inequalities for Markov chains (e.g., a' la MCDiarmid) that could perhaps provide alternative ways of estimating the mixing-relevant quantities considered here. One classic is

## Novelty and related work:

Samson, P.-M. (2000). Concentration of measure inequalities for Markov chains and  $\Phi$ -mixing processes. Ann. Probab. 28, pp. 416–461

and a more recent one is

D. Paulin, Concentration inequalities for Markov chains by Marton couplings and spectral methods. arxiv, 2015.

How are the authors positioning themselves compared to this



	(uncited!) literature ?
Soundness:	All main claimes are well suppoted.
Summary of review:	This paper proposes and analyzes a neat way of estimating the mixing time of a reversible (finite state) Markov chain in a fully data-dependent way by means of a single path dynamic of the chain. Upper and lower bounds are provided on the spectral gap and the minimum probability of the stationary distribution, the two quantities which are known to determine mixing rate. Neat, well-motivated, and well-presented piece of work.
Significance:	5: (excellent)
Novelty and related work:	4: (good)
Soundness:	5: (excellent)
Presentation:	5: (excellent)
Overall evaluation:	<b>7</b> : (accept: good paper)
Reviewer's confidence:	<b>2</b> : (medium)

Copyright © 2002–2015 EasyChair