

Department of Mathematics and Statistics

Preprint MPS-2014-11

16 April 2014

Noisy Monte Carlo: Convergence of Markov chains with approximate transition kernels

by

P.Alquier, N. FrielR.G. Everittand A. Boland



- 1. large datasets: the sample size too large. This situation is very common nowadays as huge databases can be stored at no cost. For example: in genomics the cost of sequencing has fallen by a factor ôfin past decade and a half. This has led to the wide availability of sequence data the recently announced Personal Genome Project UK aims to sequence and genomes, each consisting of 10° bases.
- 2. high-dimensional parameter spaces: the sample size the reasonable, but the number of variables is too large. For example: data assimilation in numerical weather

the target distribution. We then study the special cases of a noisy version of the Exchange algorithm (Murrayet al. (2006)), and discretized Langevin Monte Carlo in Section 3. For these noisy algorithms we prove that the total variation distance decreases with the number of iterations N, of the randomisation step in the noisy algorithm, and nd a bound on this distance in terms M. We study in detail an application to intractable likelihood problems in Section 4.

2 Noisy MCMC algorithms

It turns out that a useful answer to this question is given by the study of the stability of Markov chains. There have been a long history of research on this topic, we refer the reader to the monograph by Kartashov (1996) and the references therein. Here, we will focus on a more recent method due to Mitrophanov (2005). In order to measure the distance between P and P recall the de nition of the total variation measure between two kernels:

$$kP \quad P^{k} := \sup_{2} k_{k}$$

$$9 N_0 2 N; 0 < < 1; L > 0; 8N N_0;$$

Z
V() $\hat{P}_N(_0; d)$ V(_0) + L:

$$\mathbf{k} \mathbf{P}_{N} \mathbf{P} \mathbf{k} \mathbf{P}_{N/2} \mathbf{0}$$

Then there exists an N₁ 2 N such that any \vec{P}_N , for N N₁, is geometrically ergodic with limiting distribution $_N$ and k $_N$ k $\stackrel{l}{\underset{N'=1}{\longrightarrow}}$ 0.

(We refer the reader to (Meyn and Tweedie 1993) for the de nitiok dd_{V} theorem). Note that, in contrast to the previous result, we don't know explicitly the rate of convergence of the distance between P_{N} when N is xed. However it is possible to get an estimate of this rate (see Corollary 1 page 189 in (Ferre, Herve and Ledoux 2013)) under stronger assumptions.

2.1 Noisy Metropolis-Hastings

The Metropolis-Hastings (M-H) algorithm, sequentially draws candidate observations from a distribution, conditional only upon the last observation, thus inducing a Markov chain. The M-H algorithm is based upon the observation that a Markov chain with transition density P(;) and exhibiting detailed balance for

(jy)P(;) = (jy)P(;);

Algorithm 2 Noisy Metropolis-Hastings algorithm

for n = 0 tol do Draw $^{\theta}$ h(j _n) Draw y $^{\theta}$ F $_{0}$ () Set $_{n+1} = ^{\theta}$ with probability min($1^{(\theta)}; _{n}; y^{\theta}$)) Otherwise, set_{n+1} = $_{n}$. end for

Note that (0 ; ; y) can be thought of as a randomised versio versio versio

Obviously, we expect that is chosen in such a way that 1 and so in this case, $k_{\ 0}P^n - {}^P$

that the Langevin algorithm produces a Markov chain and Redethote the corresponding transition kernel. Note that, we generally don't h(iy) = (jy) nor $_{0}P ! (jy)$, however, under some assumption P ! for some close to when is small enough, we discuss this in more detail below.

In practice, it is often the case that g(n) cannot be computed. Here again, a natural idea is to replace log (n) by an approximation or an estimate $\log(n)$, possibly using a randomization streep F_n . This yields what we term a noisy Langevin algorithm.

Algorithm 4Noisy Langevin algorithmfor n = 0 tol do
Drawy $_{n}$ $F_{n}().$ Set $_{n+1} = _{n} + _{\overline{2}} i^{by_{n}} \log (_{n}jy) + C$ N(Q,):end for

Note that a similar algorithm has been proposed in (Welling and Teh 2011; Ahn, Korattikara and Welling 2012) in the context of big data situations, where the gradient of the logarithm of the target distribution is estimated using mini-batches of the data.

We let P denote the corresponding transition kernel arising from Algorithm 4. We now prove that the Stochastic gradient Langevin algorithm, (Algorithm 4), will converge to the discrete-time Langevin di usion with transition kernel resulting from Algorithm 3.

2.3 Towards theoretical guarantees for the noisy Langevin algorithm

In this case, the approximation guarantees are not as clear as they are for the noisy Metropolis-Hastings algorithm. To begin, there are two levels of approximation:

the transition kernel targets a distribution that might be far away from jy).

Moreover, one does not simulate at each step $\mathsf{P}\mathsf{frbout}$ rather from $\hat{\mathsf{P}}$.

The rst point requires one to control the distance betweed (jy). Such an analysis is possible. Here we refer the reader to Proposition 1 in (Dalalyan and Tsybakov 2012) and also to Roberts and Stramer (Roberts and Stramer 2002) for di erent discretization schemes. It is possible to contrko \mathbf{P} k as Lemma 2.4 illustrates.

where

Lemma 2.4

$$= \mathsf{E}_{y_n} \mathsf{F}_n \exp \frac{1}{2} - \frac{1}{2}(\mathsf{r} \log (n) + \mathsf{r}^{y_n} \log (n))^2 = 1$$

The paper by Roberts and Tweedie (1996a) contains a complete study of the chain generated by P. The problem is that it is not uniformly ergodic. So Theorem 2.1 is not the appropriate tool in this situation. However, in some situations, this chain is geometrically ergodic, and in this instance we can use Theorem 2.2 instead (moreover, note that Roberts and Tweedie (1996a) provide the function/ used in the Theorem). We provide an example of such an application in Section 3 below.

2.4 Connection with the pseudo-marginal approach

There is a clear connection between this paper and the pseudo-marginal approaches described in (Beaumont 2003) and (Andrieu and Roberts 2009). In both cases a noisy acceptance probability is considered, but in pseudo-marginal approaches this is a consequence of using an estimate of the desired target distribution at each, rather than the true value. Before proceeding further, we make precise some of the terminology used in (Beaumont 2003) and (Andrieu and Roberts 2009). These papers describe two alternative algorithms, the \Monte Carlo within Metropolis" (MCWM) approach, and \grouped independence MH" (GIMH). In both cases an unbiased importance sampling estimator, is used in place of the desired target target , however the overall algorithms proceed slightly di erently. The i(+1)th iteration of the MCWM algorithm is shown in algorithm 5.

Algorithm 5 MCWM

for n = 0 to I do Draw ⁰ $h(:j_n)$.

Draw z^0 G(:j ⁰), z G(:j), where G is an importance proposal and z^0 and z are random vectors of sizeN.

Calculate the acceptance probability, (n; 0), where b_z^N and $b_{z^0}^N$ denote the 18 11.9552 Tf 5.03 0 Tc

the auxiliary variables. The same argument holds when using any unbiased estimator of the target. As regards our focus in this paper, GIMH is something of a special case, and our framework has more in common with MCWM. We note that despite its exactness, there is no particular reason for estimators from GIMH to be more statistically e cient than those from MCWM.

where () denotes the prior distribution for For example, a naive application of the Metropolis-Hastings algorithm when proposing to move_i from $h(j_i)$ results in the acceptance probability,

$$(\ ^{\theta}; \) = \min \ 1; \frac{q \circ (y) \ (\ ^{\theta})h(j \ ^{\theta})}{q \ (y) \ (\)h(\ ^{\theta}j \)} \quad \frac{Z()}{Z(\ ^{\theta})} ;$$

$$(4)$$

depending on the intractable $ra_{Z(\ell)}^{Z(\ell)}$.

One method to overcome this computational bottleneck is to use an approximation of the likelihood (yj). A composite likelihood approximation of the true likelihood, such as that of (Besag 1974), is most commonly used. This approximation consists of a product of easily normalised full-conditional distributions. The most basic composite likelihood is the pseudo likelihood which comprised of the product of full-conditional distributions, of each

f (yj)
$$\bigvee_{i=1}^{Y^{i/j}}$$
 f (y_ijy _i;):

However this approximation of the true likelihood can give unreliable estima(figisebf and Pettitt 2004), (Friet al 2009).

3.2 Exchange Algorithm

A more sophisticated approach is to use the Exchange algorithm. **Istuat:** (2006) extended the work of M llet al. (2006) to allow inference on doubly intractable distributions using the exchange algorithm. The algorithm samples from an augmented distribution

$$(^{\theta}; y^{\theta}; jy) / f(yj) ()h(^{\theta}j)f(y^{\theta}j^{\theta})$$

whose marginal distribution fds the posterior of interest. Here the auxiliary distribution $f(y^0 j^{-n})$ is the same likelihood model in whichs de ned. By sampling from this augmented distribution, the acceptance formula simpli es, as can be seen in algorithm 6, where the normalising constants arising from the likelihood and auxiliary likelihood cancel. One di culty of implementing the exchange algorithm is the requirement to $g^{\pm}m\beta l(ej^{-n})$, perfect sampling (Propp and Wilson 1996) is often possible for Markov random eld models. However when the exchange algorithm is used with MRFs the resultant chains may not mix well. For example, Caimo and Friel (2011) used adaptive direction sampling (Gilks, Roberts and George 1994) to improve the mixing of the exchange algorithm when used with ERGM models.

Murrayet al. (2006) proposed the following interpretation of the exchange algorithm. If we compare the acceptance ratios in the M-H and Exchange algorithm, the only di erence is that the ratio of the normalising constants in the M-H acceptance $\operatorname{pr}\mathbb{Z}(a)\oplus\mathbb{Z}(y^{\ell})$ is replaced by $(y^{\ell})=q_{0}(y^{\ell})$ in the exchange probability. This ratio of un-normalised likelihoods

```
Algorithm 7 Noisy Exchange algorithm
```

```
for n = 0 tol do

Draw ^{\emptyset} h(j _{n}):

for i = 1 to N do

Drawy_{i}^{\emptyset} f (j ^{\emptyset}):

end for

De ne y _{\circ} = f y_{1}^{\emptyset}; :::; y_{N}^{\emptyset}g

Set _{n+1} = ^{\emptyset} with probability min(1^(^{\emptyset}; _{n}; y _{\circ})), where

^{(\emptyset}; _{n}; y _{\circ}) = \frac{q_{\circ}(y) (^{\emptyset})h(_{n}j_{-n})}{q_{_{n}}(y) (_{-n})h(_{-j})}\frac{1}{N}\frac{X^{V}}{_{i=1}}\frac{q_{_{\circ}}(y_{i}^{\emptyset})}{q_{_{\circ}}(y_{i}^{\emptyset})}.

Otherwise, set<sub>n+1</sub> = _{n}.

end for
```

(A3) for any and l in ,

$$\operatorname{Var}_{y^0 f(y^0 j^0)} \frac{q_n(y^0)}{q_0(y^0)} < +1$$
:

Note that whe(A1) or (A2) is satisled, we necessarily have that is a bounded set, in this case, we put = $\sup_2 k k$. This also means that $O \exp(TS)$, $q(y) \exp(TS)$ for any and S, we then put := $\exp(TS)$. Also, note that this immediately implies Assumption(A3) because in this case, $\operatorname{Var}_{f(y^0)}(q_n(y^0) = q_0(y^0))$ K

Note that Liang and Jin (2011) presents a similar algorithm to that above. However in contrast to Lemma 3.1, the results in (Liang and Jin 2011) do not explicitly provide a rate of approximation with respec NtoLemma 2.2, page 9 in (Liang and Jin 2011) only states that there exists has large enough to reach arbitrarily small accuracy.

3.4 Noisy Langevin algorithm for Gibbs random elds

The discrete-time Langevin approximation (3) is unavailable for Gibbs random elds since the gradient of the log postemolog (,jy) is analytically intractable, in general. However Algorithm 4 can be used using a Monte Carlo estimate of the gradient, as follows.

$$log((jy)) = {}^{T}s(y) \quad log(z())) + log() \quad log((y))$$

r log((jy)) = s(y) $\frac{z^{\ell}()}{z()} + r \log ()$
= s(y) $\frac{p^{s(y)}[exp^{T}s(y)]}{exp(^{T}s(y))} + r \log ()$
= s(y) $E_{yj} [s(y)] + r \log ()$ (7)

In practice, $\mathbf{E}_{y^0} \in [\mathbf{s}(y^{\ell})]$ is usually not known - an exact evaluation of this quantity would require an evaluation $\overline{\mathbf{af}}($). However, it is possible to estimate it through Monte-Carlo pimulations. If we simulate = $(\mathbf{y}_1^{\ell}; ...; \mathbf{y}_n^{\ell}) = \mathbf{f}(\mathbf{i}; \mathbf{j})$, then \mathbf{E}_{yj} [$\mathbf{s}(y)$] can be estimated using $\prod_{i=1}^{n} \mathbf{s}(\mathbf{y}_i^{\ell}) = \mathbf{n}$. This gives an estimate of the gradient from (7).

$$b^{y} \log (jy) = s(y) - \frac{1}{N} \sum_{i}^{X^{y}} s(y_{i}^{\ell}) + r \log ()$$

In turn this yield the following noisy discretized Langevin algorithm.

Algorithm 8 Noisy discretized Langevin algorithm for Gibbs random elds

```
for n = 0 tol do

for i = 1 to N do

Drawy<sup>i</sup> f (j<sub>n</sub>):

end for

De ne y<sub>n</sub> = f y<sup>l</sup><sub>1</sub>;...; y<sup>l</sup><sub>N</sub>g,

Calculate<sup>b y<sub>n</sub></sup> log (<sub>n</sub>jy) = r log (<sub>n</sub>) + s(y) \frac{1}{N} \prod_{i=1}^{P} s(y^{l}_{i}):

Set

n+1 = n + \frac{1}{2} \prod_{i=1}^{P} \log((_njy) + _n); where n are i.i.d. N (Q, ):

end for
```

Remark that in this case, the bound in Lemma 2.4 can be evaluated.

Lemma 3.3 As soon as N > $4kS^2k$ k², the in Lemma 2.4 is nite with

$$= \exp - \frac{k \log(N)}{4S^2k k^2N} = 1 + \frac{4k^p - Sk k}{N} = N! + \frac{4k \log \frac{N}{k}}{N}$$

(where k = supf k xk; kxk = 1g).

We conclude by an application of Theorem 2.2 that allows to assess the convergence of this scheme wheth ! 1 when the parameter is real.

Theorem 3.4 Assume that 2 R and the prior is Gaussian N (Q, s^2

Algorithm 9

Algorithm 10 noisy MALA-exchange

Initialise; set , for i = 1 to N do Drawy_i $f(j_0)$: end for De ne $y_0 = f y_1; \ldots; y_N g_i$ Calculater $\mathbf{b}^{y_{0}} \log (_{0}jy) = \mathbf{r} \log (_{0}) + \mathbf{s}(y) \frac{1}{N} \mathbf{P}_{i=1}^{N} \mathbf{s}(y_{i})$: for n = 0 tol do Draw $^{\ell} = n + \frac{1}{2} \mathbf{b}^{y_n} \log (n_j \mathbf{y}) +$ N(Q,). for i = 1 to N do Draw y_i^{ℓ} f (j^{ℓ}): end for de ne y $_{0} = f y_{1}^{\emptyset}; \ldots; y_{N}^{\emptyset} g.$ Calculate^{b y} \circ log ($^{\ell}jy$) = r log ($^{\ell}$) + s(y) $\frac{1}{N} \stackrel{\mathsf{P}}{\underset{i=1}{\overset{N}{\longrightarrow}}} s(y_i^{\ell})$: Set $_{n+1} = {}^{\ell}$ and $y_{n+1} = y_0$ with probability min(1^(l); $_n; y_n$)) where $(\hat{y}_{n}; y_{n}; y_{n}) = \frac{q_{0}(y) (\hat{y}_{n})h(\hat{y}_{n}; y_{n}^{\ell})}{q_{n}(y) (\hat{y}_{n})h(\hat{y}_{n}; y_{n}^{\ell})} \frac{1}{N} \frac{X^{V}}{q_{n}\overline{n}} q_{n}\overline{n}$ %87is9701 Tf 5.]TJ/2F34 7.9701 Tf

4.1 Ising study

The Ising model is de ned on a rectangular lattice or grid. It is used to model the spatial distribution of binary variables, taking values and 1. The joint density of the Ising model can be written as

$$f(yj) = \frac{1}{Z()} \exp \left(\begin{array}{c} X^{j} X \\ y_{i}y_{j} \end{array} \right)$$

where j denotes that and j are neighbours an $\mathbf{Z}() = \mathbf{P}_{y} \exp^{n} \mathbf{P}_{M} \exp^{n} \mathbf{P}_{j=1} \mathbf{P}_{j} \mathbf{y}_{j} \mathbf{y}_{j}$. The normalising constact () is rarely available analytically since this relies on taking the summation over all di erent possible realisations of the lattice. For a lattice modes this equates to $\frac{M(M-1)}{2^{2}}$ di erent possible lattice formations.

For our study, we simulated 20 grids of size166 This size lattice is su ciently small enough such that the normalising cons26(n)tcan be calculated exactly (36.5 minutes for each graph) using a recursive forward-backward algorithm (Reeves and Pettitt 2004; Friel and Rue 2007), giving a gold standard with which to compare the other algorithms. This is done by calculating the exact density over a ne gridvolues, f₁; , g over the interval [04;08], which cover the e ective range of values that take. We normalise(,jy) by numerically integrating over the un-normalised density.

$$^{(y)} = \frac{X'}{\sum_{i=2}^{j}} \frac{(i - i - 1)}{2} \frac{q_i(y)}{Z(i)} (i) + \frac{q_{i-1}(y)}{Z(i-1)} (i - 1); \qquad (8)$$

yielding

$$(_{j}jy) = \frac{q_{i}(y)}{Z(_{j})} \frac{(_{j})}{(_{j})}$$

Each of the algorithms was run for 30 seconds on each of the 20 datasets, at each iteration the auxiliary step to draw



Figure 1: Boxplot of the bias estimate for 20 datasets corresponding to the exchange, importance sampling exchange, Langevin and MALA algorithms.

Figure 1 shows the bias of the posterior means for each of the algorithms. We see that both

Figure 2: Estimated posterior densities corresponding to the exact and noisy algorithms corresponding to one of the datasets used in the Ising simulation study.

4.2 ERGM study

Here we explore how our algorithms may be applied to the exponential random graph model (ERGM) (Robins **et al**

4.2.1 The Florentine Business dataset

Here, we consider a simple 16 node undirected graph: the Florentine family business graph. This concerns the business relations between some Florentine families in around 1430. The network is displayed in Figure 3. We propose to estimate the following 2-dimensional model.

$$f(yj) = \frac{1}{Z()} \exp(_{1}s_{1}(y) + _{2}s_{2}(y));$$

where $s_1(y)$ is the number of edges in the graph $s_2(y)$ is the number of two-stars.



Before we could run the algorithms, certain parameters had to be tuned. We used a at prior N (Q,100) in all of the algorithms. The Langevin, MALA exchange and noisy MALA exchange algorithms all depend on a stepsize matrix . This matrix determines the scale of proposal values for each of the parameters. This matrix should be set up so that proposed values for accommodate the di erent scales of the posterior densitynof order to have good mixing in the algorithms we chose a which relates to the shape of the posterior density. Our approach was to aim to relate to the covariance of the posterior density. To do this, we equated to an estimate of the inverse of the second derivative of the log posterior at the algorithm a posteriori estimate . As the true value of the MAP is unknown, we used a Robbins-Monro algorithm (Robbins and Monro 1951) to estimate this. The Robbins-Monro algorithm takes steps in the direction of the slope of the distribution. It is very similar to Algorithm 8 except without the added noise and follows the stochastic process

where
$$\int_{i=0}^{n+1} = n + n e^{j y_n} \log (n j y);$$

 $X^V \qquad X^V \qquad X^V \qquad 2 < 1 :$

The values of decrease over time and once the di erence between successive values of this process is less than a speci ed tolerance level, the algorithm is deemed to have converged to the MAP. The second derivative of the log posterior is derived by di erentiating (7) yielding

r²log (j

MALA exchange but not in the Langevin algorithm. Since our Noisy Langevin algorithm approximates Langevin di usion we are approximating an approximation. There are two levels of approximations which leaves more room for error.

Edge 2-star

Method

Figure 5: Chains, density plot and ACF plot for the 2-star statistic.

4.2.2 The Molecule dataset

The Molecule dataset is a 20 node graph, shown in Figure 6. We consider a four parameter model which includes the number of edges in the graph, the number of two-stars, the number of three-stars and the number of triangles.

$$f(yj) = \frac{1}{Z()} \exp(_{1}s_{1}(y) + _{2}s_{2}(y) + _{3}s_{3}(y) + _{4}s_{4}(y))$$

The parameter was chosen in a similar fashion to the Florentine business example. The Robbins-Monro algorithm was run for 20,000 iterations to nd an estimate of the MAP, 4,000 graphs were then simulated at the estimated MAP and these were used to calculate an estimate of the second derivative using Equation (9). The matrix was theofanth estimated



Figure 6: Molecule network

The BERGM algorithm of (Caimo and Friel 2011) was again used as a \ground truth". This algorithm was run for a large number of iterations equating to 4 hours of CPU time. This gave us accurate estimates against which to compare the various algorithms. The ve algorithms were each run for 100 seconds of CPU time. Table 2 shows the posterior mean and standard deviations of each of the four parameters for each of the algorithms. The results for the Molecule dataset model are similar to the Florentine business dataset model. In Table 2 we see that the noisy exchange algorithm improved on the standard exchange algorithm. The MALA exchange improved on noisy Langevin and the Noisy MALA improved on the MALA exchange.

Figure 7 and Figure 8 show the densities and the autocorrelation plots of the algorithms. The autocorrelation plots show that the noisy algorithms had less correlation than the exchange algorithm. The densities show that again the algorithms, when run on the Molecule model, performed in the same manner as the Florentine model. The algorithms with the exception of the noisy Langevin algorithm estimated the mode well but underestimated the standard deviation. The noisy Langevin algorithm did not estimate the mean or standard deviations well.

	Edge 2-star		3-Star		Triangle			
Method	Mean	SD N	/lean SD	Mea	an SD	Mean	SD	
BERGM	2.647	2.754	-1.069	0.953	-0.021	0.483	1.787	0.646
Exchange	1.889	2.142	-0.797	0.744	-0.138	0.385	1.593	0.519
Noisy Exch	1.927	2.444	-0.757	0.823	-0.176	0.422	1.543	0.53
Noisy Lang	1.679	3.65	-0.509	1.429	-0.466	0.787	1.633	0.573
MALA Exch	2.391	2.095	-0.938	0.795	-0.113	0.451	1.454	0.598
Noisy MALA Exch	2.731	2.749	-1.054	0.886	-0.041	0.417	1.519	0.492

Table 2: Posterior means and standard deviations.

Figure 7: Density plots of the 4 parameters for the molecule example.

result to hold for noisy MCMC algorithms, in which case the e ect of this additional variance on top of the aforementioned bias should be a consideration when employing noisy MCMC.

A further area for future work lies in relaxing the requirement for the ideal non-noisy chain to be uniformly ergodic. This property does not hold in many cases: the results in this paper are intended as the rst steps towards future work that would obtain results that hold more generally.

Acknowledgements

The Insight Centre for Data Analytics is supported by Science Foundation Ireland under Grant Number SFI/12/RC/2289. Nial Friel's research was also supported by an Science Foundation Ireland grant: 12/IP/1424.

References

- Ahn, S., A. Korattikara and M. Welling (2012), Bayesian Posterior Sampling via Stochastic Gradient Fisher Scoring. IProceedings of the 29th International Conference on Machine Learning
- Andrieu, C. and G. Roberts (2009), The pseudo-marginal approach for e cient Monte-Carlo computation **The Annals of Statistics 37**(2), 697{725
- Andrieu, C. and M. Vihola (2012), Convergence properties of pseudo-marginal Markov

- Friel, N. and A. N. Pettitt (2004), Likelihood estimation and inference for the autologistic model. Journal of Computional and Graphical Statistics 3, 232{246
- Friel, N., A. N. Pettitt, R. Reeves and E. Wit (2009), Bayesian inference in hidden Markov random elds for binary data de ned on large latti**deux**nal of Computational and Graphical Statistics 18, 243{261
- Friel, N. and H. Rue (2007), Recursive computing and simulation-free inference for general factorizable models.

$$ZZ = h = (d^{0}) dt dy^{0}h(tj)F_{t}(y^{0}) min(1;^{(};t;y^{0})) min(1;(;t)) Z = h + dy^{0}F_{0}(y^{0}) h(^{0}j) min(1;(;^{0})) h(^{0}j) min(1;^{(}; d)) d$$

Now, note that

$$Z = \frac{1}{p \cdot \frac{1}{2}} \exp \frac{ktk^2}{2} + 1 = \#$$

$$\exp \frac{t^7 \cdot \frac{1}{2}(r \log () - r^{^{^{^{^{^{0}}}}}\log ())}}{2} + \frac{1}{8}k^{-\frac{1}{2}}(r \log () - r^{^{^{^{^{^{0}}}}}\log ())k^2 + dt$$

$$= E + 1 \exp a^7 X + \frac{kak^2}{2}$$

where X N (0,1) and $a = \frac{1}{2} [r \log () r^{5}y^{0} \log ()]=2$. Then:

E1 exp
$$a^T X$$
 $\frac{kak^2}{2} = exp$ $\frac{kak^2}{2}$ E exp $a^T X$ exp $\frac{kak^2}{2}$

$$= exp$$
 $\frac{kak^2}{2}$ E exp $a^T X$ E exp $a^T X$

$$exp$$
 $\frac{kak^2}{2}$ P $\overline{Var[exp(TX)]}$

$$= exp$$
 $\frac{kak^2}{2}$ Q $\overline{E[exp(a^T X)]}$ E $[exp(a^T X)]^2$

$$= exp$$
 $\frac{kak^2}{2}$ P $\overline{exp(kak^2)}$ exp (kak^2)

$$= \frac{P}{exp(kak^2)}$$
 1:

So nally,

kP P

$$p \frac{1}{\overline{N}} \frac{h(j \circ) (\circ) q \circ (y)}{h(\circ) (\circ) q \circ (y)} s \frac{s}{Var_{y_1^0 \circ f(y_1^0) \circ} \frac{q_n(y_1^0)}{q \circ (y_1^0)}} :$$

Proof of Theorem 3.2. Under the assumptions of Theorem 3.2, note that (4) leads to

$$(_{n}; ^{0}) = \frac{(^{0})q_{0}(y)Z(_{n})}{(_{n})q_{n}(y)Z(^{0})} \frac{h(_{n}j ^{0})}{h(^{0})} - \frac{1}{c^{2}c_{h}^{2}K^{4}} :$$
 (10)

Let us consider any measurable subset of and 2. We have

$$P(;B) = \begin{bmatrix} Z & Z \\ (d^{0}) & 1 & dth(tj)min(1; (;t)) \\ B & Z \\ + & d^{0}h \end{pmatrix} 1$$

with

$$C = c^2 c_h^2 K^4 + \frac{C}{1}$$

with $= \frac{I \underset{log(1=C)}{0} m}{\frac{1}{\log(1)}}$. Proof of Lemma 3.3. Note that

r log ()
$$r^{\Lambda_{X^0}} = \frac{1}{N} \sum_{i=1}^{X^V} s(y_i^0) = E_{y^0 f} [s(y^0)]$$

So we have to nd an upper bound, uniformly overor

$$\mathsf{D} := \mathsf{E}_{y^0 \ F_n} : \exp^4 \frac{2}{2} \quad \frac{1}{2} \quad \frac{1}{\mathsf{N}} \overset{\mathsf{X}^{\mathsf{V}}}{\underset{i=1}{\mathsf{X}^{\mathsf{V}}}} \mathsf{s}(y^{\emptyset}_i) \quad \mathsf{E}_{y^0 \ f} \ [\mathsf{s}(y^{\emptyset})] \quad 5 \quad 1; :$$

Let us put $V := \frac{1}{N} \bigvee_{i=1}^{N} V^{(i)} := \frac{1}{N} \bigvee_{i=1}^{N} \bigvee_{i=1}^{\frac{1}{2}} f_{\mathbf{s}}(\mathbf{y}_{i}^{\ell}) \quad \mathbf{E}_{y^{0} \ f} \ [\mathbf{s}(\mathbf{y}^{\ell})]\mathbf{g}$ and denote $\mathbf{W}_{j} \ (\mathbf{j} = 1; \dots; \mathbf{k})$

exp k

and soP is geometrically ergodic with functionWe calculate

$$Z = \frac{2}{V()P^{A}(_{0};d)} = E_{y^{0}}4\frac{p}{\frac{1}{2}} = \frac{2}{R}V()\exp@ - \frac{0}{2}r^{y^{0}}\log(_{0}jy) = \frac{1}{A}d^{5}$$
$$= E_{y^{0}}\frac{p}{\frac{1}{2}} = \frac{2}{R}V + \frac{2}{2}(r^{y^{0}}\log(_{0}jy) = \frac{1}{A}d^{5})$$