

# Approximate spectral gaps for Markov chains mixing times in high dimensions\*

Yves F. Atchadé<sup>†</sup>

**Abstract.** This paper introduces a concept of approximate spectral gap to analyze the mixing time of reversible Markov Chain Monte Carlo (MCMC) algorithms for which the usual spectral gap is degenerate or almost degenerate. We use the idea to analyze a MCMC algorithm to sample from mixtures of densities. As an application we study the mixing time of a Gibbs sampler for variable selection in linear regression models. We show that properly tuned, the algorithm has a mixing time that grows at most polynomially with the dimension. Our results also suggest that the mixing time improves when the posterior distribution contracts towards the true model and the initial distribution is well-chosen.

**Key words.** Markov Chain Monte Carlo algorithms, Markov chains mixing times, Spectral gaps, Canonical paths, MCMC for mixtures of densities, High-dimensional linear regression models

**AMS subject classifications.** 60J05, 65C05, 65C60

**1. Introduction.** Understanding the type of problems for which fast Markov Chain Monte Carlo (MCMC) sampling is possible is a question of fundamental interest. The study of the size of the spectral gap is a widely used approach to gain insight into the behavior of MCMC algorithms. However this technique may be inappropriate when dealing with distributions with small isolated local modes. To be more precise, let  $\pi$  be some probability measure of interest on some measure space  $\mathcal{X}$ , and let  $K$  be a Markov kernel with invariant distribution  $\pi$ . For the purpose of sampling from  $\pi$  using  $K$ , one can represent an isolated local mode (to which  $K$  is sensitive) as a subset  $A$  such that  $K(x, \mathcal{X} \setminus A)$  is small compared to  $\pi(\mathcal{X} \setminus A)$  for all  $x \in A$ . In this case,  $K$  will have a small conductance, and a small spectral gap. Note however that if  $\pi(A)$  is also small (that is we are dealing with a small isolated mode  $A$ ), then, since

$$\int_{\mathcal{X} \setminus A} \pi(dx) K(x, A) = \int_A \pi(dx) K(x, \mathcal{X} \setminus A),$$

we see that the set  $A$  will be typically hard to reach in the first place. Hence, any finite-length Markov chain  $\{X_0, \dots, X_n\}$  say, with transition kernel  $K$  and initialized in  $\mathcal{X} \setminus A$  is unlikely to visit  $A$ . Nevertheless, and since  $\pi(A)$  is small,  $X_n$  may still be a good approximate sample from  $\pi$  for large  $n$ . This implies that the poor mixing time predicted by the standard spectral gap may markedly differ from the actual behavior of these finite-length chains. Motivated by this problem, and building on the  $s$ -conductance of L. Lovasz and M. Simonovits ([Lovász and Simonovits(1993)]), we develop an idea of approximate spectral gap (that we call  $\zeta$ -spectral gap, for some  $\zeta \in [0, 1)$ ) which allows us to measure the mixing time of a Markov chain while discounting the ill-effect of overly small (and potentially problematic) sets.

\*

**Funding:** This work is partially supported by the NSF grant DMS1513040.

<sup>†</sup>Boston University, Department of Mathematics and Statistics, 111 Cummington Mall, Boston, 02215, MA, United States ([atchade@bu.edu](mailto:atchade@bu.edu), <http://math.bu.edu/people/atchade/>).

Mixtures are good examples of probability distributions with isolated local modes. We use the idea to analyze a class of MCMC algorithms to sample from mixtures of densities. Much is known on the computational complexity of various MCMC algorithms for log-concave densities (see e.g. [Lovász and Simonovits(1993), Frieze et al.(1994), Lovász(1999), Lovász and Vempala(2007)], and [Dwivedi et al.(2018)] and the references therein). However these results cannot be directly applied to mixtures, since a mixture of log-concave densities is not log-concave in general. By augmenting the variable of interest to include the mixing variable, a Gibbs sampler can be used to sample from a mixture. A very nice lower bound on the spectral gap of such Gibbs samplers is developed in [Madras and Randall(2002)]. We re-examine [Madras and Randall(2002)]'s argument using the concept of  $\zeta$ -spectral gap, leading to Theorem 3.1 that gives potentially better dependence on the dimension.

Our initial motivation into this work is in large-scale Bayesian variable selection problems. The Bayesian posterior distributions that arise from these problems are typically mixtures of log-concave densities with very large numbers of components, and the aforementioned Gibbs sampler is commonly used for sampling (see e.g. [George and McCulloch(1997), Narisetty and He(2014)]). We show that when properly tuned, the algorithm has a mixing time that grows at most polynomially with  $p$ , the number of regressors in the model (Theorem 4.2). Our result derived from the approximate spectral gap also suggests that the mixing time improves when a good initial distribution is used, provided that posterior contraction towards the true model holds (Theorem 4.3).

The paper is organized as follows. We develop the concept of  $\zeta$ -spectral gap in Section 2. The main result there is Lemma 2.1. In Section 3 we study the mixing time of mixtures of Markov kernels, and derive (Theorem 3.1) a generalization of Theorem 1.2 of [Madras and Randall(2002)]. We put these two results together to analysis the linear regression model in Section 4. Some numerical simulations are detailed in Section 4.1.

**2. Approximate spectral gaps for Markov chains.** Let  $\pi$  be a probability measure on some Polish space  $(\mathcal{X}, \mathcal{B})$  (where  $\mathcal{B}$  is its Borel sigma-algebra), equipped with a reference sigma-finite measure denoted  $dx$ . In the applications that we have in mind,  $\mathcal{X}$  is the Euclidean space  $\mathbb{R}^p$  equipped with its Lebesgue measure. We assume that  $\pi$  is absolutely continuous with respect to  $dx$ , and we will abuse notation and use  $\pi$  to denote both  $\pi$  and its density:  $\pi(dx) = \pi(x)dx$ . We let  $L^2(\pi)$  denote the Hilbert space of all real-valued square-integrable (wrt  $\pi$ ) functions on  $\mathcal{X}$ , equipped with the inner product  $\langle f, g \rangle_\pi \stackrel{\text{def}}{=} \int_{\mathcal{X}} f(x)g(x)\pi(dx)$  with associated norm  $\|\cdot\|_{2,\pi}$ . More generally, for  $s \geq 1$ , we set  $\|f\|_{s,\pi} \stackrel{\text{def}}{=} (\int_{\mathcal{X}} |f(x)|^s \pi(dx))^{1/s}$ . For  $s = +\infty$ ,  $\|f\|_{s,\pi}$  is defined as the essential supremum of  $|f|$  with respect to  $\pi$ . If  $P$  is a Markov kernel on  $\mathcal{X}$ , and  $n \geq 1$  an integer,  $P^n$  denotes the  $n$ -th iterate of  $P$ , defined recursively as  $P^n(x, A) \stackrel{\text{def}}{=} \int_{\mathcal{X}} P^{n-1}(x, dz)P(z, A)$ ,  $x \in \mathcal{X}$ ,  $A$  measurable. If  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a measurable function, then  $Pf : \mathcal{X} \rightarrow \mathbb{R}$  is the function defined as  $Pf(x) \stackrel{\text{def}}{=} \int_{\mathcal{X}} P(x, dz)f(z)$ ,  $x \in \mathcal{X}$ , assuming that the integral is well defined. And if  $\mu$  is a probability measure on  $\mathcal{X}$ , then  $\mu P$  is the probability on  $\mathcal{X}$  defined as  $\mu P(A) \stackrel{\text{def}}{=} \int_{\mathcal{X}} \mu(dz)P(z, A)$ ,  $A \in \mathcal{B}$ . The total variation distance between two probability measures  $\mu, \nu$  is defined as

$$\|\mu - \nu\|_{\text{tv}} \stackrel{\text{def}}{=} 2 \sup_{A \in \mathcal{B}} (\mu(A) - \nu(A)).$$

76 Let  $K$  be a Markov kernel on  $\mathcal{X}$  that is reversible with respect to  $\pi$ . That is for all  
 77  $A, B \in \mathcal{B}$ ,

$$78 \quad \int_A \pi(dx) \int_B K(x, dy) = \int_B \pi(dx) \int_A K(x, dy).$$

79 We will also assume throughout that  $K$  is lazy in the sense that  $K(x, \{x\}) \geq \frac{1}{2}$ . The concept  
 80 of spectral gap and the related Poincaré's inequalities are commonly used to quantify Markov  
 81 chains' mixing times. For  $f \in L^2(\pi)$ , we set  $\pi(f) \stackrel{\text{def}}{=} \int_{\mathcal{X}} f(x)\pi(dx)$ ,  $\text{Var}_{\pi}(f) \stackrel{\text{def}}{=} \|f - \pi(f)\|_{2,\pi}^2$ ,  
 82 and  $\mathcal{E}(f, f) \stackrel{\text{def}}{=} \frac{1}{2} \int \int (f(y) - f(x))^2 \pi(dx) K(x, dy)$ . The spectral gap of  $K$  is then defined as

$$83 \quad \text{SpecGap}(K) \stackrel{\text{def}}{=} \inf \left\{ \frac{\mathcal{E}(f, f)}{\text{Var}_{\pi}(f)}, f \in L^2(\pi), \text{ s.t. } \text{Var}_{\pi}(f) > 0 \right\}.$$

84 It is well-known (see for instance [Montenegro and Tetali(2006)] Corollary 2.15) that if  $\pi_0(dx) =$  ■  
 85  $f_0(x)\pi(dx)$ , and  $f_0 \in L^2(\pi)$ , then

$$86 \quad (2.1) \quad \|\pi_0 K^n - \pi\|_{\text{tv}}^2 \leq \text{Var}_{\pi}(f_0) (1 - \text{SpecGap}(K))^n.$$

87 Therefore, lower-bounds on the spectral gap can be used to derive upper-bounds on the mixing  
 88 time of  $K$ . In many examples, the conductance of  $K$  is easier to control than the spectral gap.  
 89 In these examples the concept of  $s$ -conductance introduced by L. Lovász and M. Simonovits  
 90 ([Lovász and Simonovits(1993)]) as a generalization of the conductance has proven very useful,  
 91 particularly in problems where a warm-start to the Markov chain is available. For  $\zeta \in [0, 1/2)$ ,  
 92 we define the  $\zeta$ -conductance of the Markov kernel  $K$  as

$$93 \quad \Phi_{\zeta}(K) \stackrel{\text{def}}{=} \inf \left\{ \frac{\int_A \pi(dx) K(x, A^c)}{(\pi(A) - \zeta)(\pi(A^c) - \zeta)}, \zeta < \pi(A) < \frac{1}{2} \right\},$$

94 where the infimum above is taken over measurable subsets of  $\mathcal{X}$ . Note that  $\Phi_0(K)$  is the  
 95 standard conductance. Plainly put,  $\Phi_{\zeta}(K)$  captures the same concept of ergodic flow as  
 96  $\Phi_0(K)$ , except that in  $\Phi_{\zeta}(K)$  we disregard sets that are either too small or too large under  $\pi$ .  
 97 It turns out that  $\Phi_{\zeta}(K)$  still controls the mixing time of  $K$  up to an additive constant that  
 98 depends on  $\zeta$  (see [Lovász and Simonovits(1993)] Corollary 1.5). One important drawback of  
 99 the  $\zeta$ -conductance is that the arguments that relate  $\Phi_{\zeta}(K)$  to the mixing time of  $K$  (Theorem  
 100 1.4 of [Lovász and Simonovits(1993)]) is rather involved, and this has limited the scope and  
 101 the usefulness of the concept. Furthermore there are some problems where direct bound on  
 102 the spectral gap instead of the conductance is easier, or yields better results.

103 Motivated by the  $\zeta$ -conductance, we introduce a similar concept of  $\zeta$ -spectral gap that  
 104 directly approximates the spectral gap. Let  $\|\cdot\|_{\star} : L^2(\pi) \rightarrow [0, \infty]$  denote a norm-like  
 105 function on  $L^2(\pi)$  with the following properties: (i)  $\|\alpha f\|_{\star} = |\alpha| \|f\|_{\star}$ , (ii) if  $\|f\|_{\star} = 0$  then  
 106  $\text{Var}_{\pi}(f) = 0$ , and (iii)

$$107 \quad (2.2) \quad \|Kf\|_{\star} \leq \|f\|_{\star}, \quad f \in L_{\star}^2(\pi),$$

108 where  $L_{\star}^2(\pi) \stackrel{\text{def}}{=} \{f \in L^2(\pi) : \|f\|_{\star} < \infty\}$ . For  $\zeta \in (0, 1)$ , we define the  $\zeta$ -spectral gap of  $K$  as

$$109 \quad (2.3) \quad \text{SpecGap}_{\zeta}(K) \stackrel{\text{def}}{=} \inf \left\{ \frac{\mathcal{E}(f, f)}{\text{Var}_{\pi}(f) - \frac{\zeta}{2}}, f \in L_{\star}^2(\pi), \text{Var}_{\pi}(f) > \zeta, \text{ and } \|f\|_{\star} = 1 \right\}.$$

110 We note that  $\text{SpecGap}_\zeta(K)$  depends on the choice of  $\|\cdot\|_\star$ . We note also that if  $\zeta = 0$  and  
 111  $\|f\|_\star = \|f\|_{2,\pi}$ , then we recover  $\text{SpecGap}_0(K) = \text{SpecGap}(K)$ . Furthermore, given  $f \in L^2(\pi)$ ,  
 112 and writing  $\bar{f} = f - \pi(f)$ , we have

$$113 \quad \frac{\mathcal{E}(f, f)}{\text{Var}_\pi(f) - \frac{\zeta}{2}} = \frac{\pi(\bar{f}^2) - \langle \bar{f}, P\bar{f} \rangle_\pi}{\pi(\bar{f}^2) - \frac{\zeta}{2}}.$$

114 By the laziness of the chain,  $\langle \bar{f}, P\bar{f} \rangle_\pi \geq \pi(\bar{f}^2)/2$ , and we deduce that  $\text{SpecGap}_\zeta(K)$  is a  
 115 quantity that always belongs to the interval  $[0, 1]$ . The idea is somewhat similar to the con-  
 116 cept of weak Poincare inequality developed for continuous-time Markov semigroups with zero  
 117 spectral gap ([Liggett(1991), Cattiaux and Guillin(2009)]). One key difference is that weak  
 118 Poincare inequalities lead to sub-geometric rates of convergence of the semigroup, whereas  
 119 the idea of  $\zeta$ -spectral gap as introduced here leads to a geometric convergence rate, plus an  
 120 additive remainder that depends on  $\zeta$ . More precisely, we have the following analog of (2.1).  
 121 The proof is similar to the proof of (2.1).

122 **Lemma 2.1.** *Suppose that  $K$  is  $\pi$ -reversible, lazy, and satisfies (2.2). Fix  $\zeta \in [0, 1)$ . Sup-*  
 123 *pose that  $\pi_0(dx) = f_0(x)\pi(dx)$  for a function  $f_0 \in L^2_\star(\pi)$ . Then for all integer  $n \geq 1$ , we*  
 124 *have*

$$125 \quad \|\pi_0 K^n - \pi\|_{\text{tv}}^2 \leq \text{Var}_\pi(K^n f_0) \leq \text{Var}_\pi(f_0) (1 - \text{SpecGap}_\zeta(K))^n + \zeta \|f_0\|_\star^2.$$

126 *Proof.* See Section 5.1. ■

127 It is also possible to control similarly the convergence to stationarity in the 1-Wasserstein  
 128 metric. Indeed, for any  $h \in L^2(\pi)$  we have

$$129 \quad (2.4) \quad |\pi_0 K^n(h) - \pi(h)| = \left| \int_{\mathcal{X}} h(x) (K^n f_0(x) - 1) \pi(dx) \right| \leq \|h\|_{2,\pi} \sqrt{\text{Var}_\pi(K^n f_0)}.$$

130 Hence, if  $\mathcal{X}$  is a metric space and  $\pi$  is such that any Lipschitz function  $h$  on  $\mathcal{X}$  belongs to  
 131  $L^2(\pi)$  (basically  $\pi$  has finite second moments), then under the assumptions of Lemma 2.1 we  
 132 have,

$$134 \quad (2.5) \quad \mathbb{W}_1(\pi_0 K^n, \pi) \stackrel{\text{def}}{=} \sup_{h: \|h\|_{\text{Lip}}=1} |\pi_0 K^n(h) - \pi(h)|$$

$$135 \quad \leq \sup_{h: \|h\|_{\text{Lip}}=1} \|h\|_{2,\pi} \sqrt{\text{Var}_\pi(f_0) (1 - \text{SpecGap}_\zeta(K))^n + \zeta \|f_0\|_\star^2},$$

137 where  $\|h\|_{\text{Lip}} \stackrel{\text{def}}{=} \sup_{x \neq y} |h(y) - h(x)|/d(y, x)$  is the Lipschitz norm of  $h$ , and where  $d$  is the  
 138 metric on  $\mathcal{X}$ .

139 **2.1. Illustration with the small local mode example.** We now illustrate how the approx-  
 140 imate spectral gap can be used with the conceptual example described in the introduction.  
 141 For that purpose, in this section we assume that  $\mathcal{X} = \mathcal{X}_0 \cup (\mathcal{X}_0^c)$  for some measurable subset  
 142  $\mathcal{X}_0$  of  $\mathcal{X}$ . We aim to capture the intuition that when  $\mathcal{X}_0^c$  is small under  $\pi$ , a Markov chain  
 143 with transition kernel  $K$  started in  $\mathcal{X}_0$  typically does not suffer from the local modes in  $\mathcal{X}_0^c$ .

144 Let  $\mathcal{B}_{\mathcal{X}_0}$  be the trace sigma-algebra of  $\mathcal{B}$  on  $\mathcal{X}_0$ . Let  $K_{\mathcal{X}_0}$  be the restriction of  $K$  on  $\mathcal{X}_0$ .  
 145 That is  $K_{\mathcal{X}_0}$  is the transition kernel on  $(\mathcal{X}_0, \mathcal{B}_{\mathcal{X}_0})$  defined as

$$146 \quad K_{\mathcal{X}_0}(x, dy) = K(x, dy) + \delta_x(dy)K(x, \mathcal{X}_0^c), \quad x \in \mathcal{X}_0.$$

147 Using the reversibility of  $K$ , it is easy to show that the invariant distribution of  $K_{\mathcal{X}_0}$  is  $\pi_{\mathcal{X}_0}$ ,  
 148 the restriction of  $\pi$  to  $\mathcal{X}_0$ , and the spectral gap of  $K_{\mathcal{X}_0}$  is given by

$$149 \quad (2.6) \quad \text{SpecGap}_{\mathcal{X}_0}(K) \stackrel{\text{def}}{=} \inf \left\{ \frac{1 \int_{\mathcal{X}_0} \int_{\mathcal{X}_0} \pi(dx) K(x, dy) (f(y) - f(x))^2}{2 \frac{1}{2} \int_{\mathcal{X}_0} \int_{\mathcal{X}_0} \pi(dx) \pi(dy) (f(y) - f(x))^2}, f : \mathcal{X} \rightarrow \mathbb{R} \right\},$$

150 where the infimum is taken over all functions  $f \in L^2_\star(\pi)$  such that

151  $\int_{\mathcal{X}_0} \int_{\mathcal{X}_0} \pi(dx) \pi(dy) (f(y) - f(x))^2 > 0$ . The next result shows that the spectral gap of  $K_{\mathcal{X}_0}$  is  
 152 a lower bound for  $\text{SpecGap}_\zeta(K)$ .

153 **Lemma 2.2.** *For  $\zeta \in (0, 1)$ , and  $\|\cdot\|_\star = \|\cdot\|_{m, \pi}$ , for some  $m \in (2, +\infty]$ , if  $\pi(\mathcal{X}_0) \geq$   
 154  $1 - \left(\frac{\zeta}{10}\right)^{1 + \frac{2}{m-2}}$  then we have*

$$155 \quad \text{SpecGap}_\zeta(K) \geq \text{SpecGap}_{\mathcal{X}_0}(K).$$

156 *Proof.* See Section 5.2. ■

157 Fix  $\zeta_0 \in (0, 1)$ . Suppose that we choose the initial distribution  $\pi_0$  such that  $\|f_0\|_{m, \pi} \leq B$ ,  
 158 for some constant  $B \geq 1$ . In that case Lemma 2.1 with  $\|\cdot\|_\star = \|\cdot\|_{m, \pi}$ , and  $\zeta = \zeta_0^2/(B^2)$  gives  
 159 for all  $n \geq 1$ ,

$$160 \quad (2.7) \quad \|\pi_0 K^n - \pi\|_{\text{tv}}^2 \leq B^2 (1 - \text{SpecGap}_\zeta(K))^n + \zeta_0^2.$$

161 Therefore, if  $\pi(\mathcal{X}_0) \geq 1 - \left(\frac{\zeta}{10}\right)^{1 + \frac{2}{m-2}}$ , by Lemma 2.2 we obtain the following bound on the  
 162 mixing time:

$$163 \quad \|\pi_0 K^N - \pi\|_{\text{tv}} \leq \sqrt{2}\zeta_0, \quad \text{for all } N \geq \frac{\log\left(\frac{B^2}{\zeta_0}\right)}{\text{SpecGap}_{\mathcal{X}_0}(K)}.$$

164 In other words the mixing time of  $K$  can indeed be controlled by the spectral gap of  $K_{\mathcal{X}_0}$ . The  
 165 condition  $\pi(\mathcal{X}_0) \geq 1 - \left(\frac{\zeta_0^2}{10B^2}\right)^{1 + \frac{2}{m-2}}$  puts a stringent constraint on the initial distribution  $\pi_0$   
 166 and on the concentration properties of  $\pi$  on  $\mathcal{X}_0$ . The successful use of the technique typically  
 167 hinges on controlling these two aspects. Further illustrations are given below.

168 **2.2. Extension to reversible Markov semigroups.** The idea can also be applied to continuous-  
 169 time Markov processes. We refer the reader to ([Bakry et al.(2013)]) for an introduction to  
 170 Markov semigroups. We consider a reversible Markov semigroup  $K = \{K_t, t \geq 0\}$ , where for  
 171 each  $t$ ,  $K_t$  is a Markov kernel on  $(\mathcal{X}, \mathcal{B})$  that is reversible with respect to  $\pi$ . Let  $G$  denote  
 172 the generator of the semi-group that we assumed well-defined on a dense subspace  $\mathcal{A}$  of  $L^2(\pi)$   
 173 that is stable under  $G$  and  $K_t$  such that for all  $t \geq 0$ ,

$$174 \quad (2.8) \quad \frac{d}{dt} K_t f = K_t G f = G K_t f, \quad f \in \mathcal{A}.$$

175 We make also the assumption that the domain  $\mathcal{A}$  contains constant functions and is equipped  
 176 with a norm  $\|\cdot\|_*$  such that  $\|f\|_* = 0$  implies that  $\text{Var}_\pi(f) = 0$ , and for all  $t \geq 0$

$$177 \quad (2.9) \quad \|K_t f\|_* \leq \|f\|_*, \quad f \in \mathcal{A}.$$

178 The Dirichlet form of  $K$  is defined as

$$179 \quad \mathcal{E}(f, f) \stackrel{\text{def}}{=} - \int_{\mathcal{X}} f(x) G f(x) \pi(dx).$$

180 For  $\zeta \in [0, 1)$ , we can define the  $\zeta$ -spectral gap of the semi-group  $K$  as

$$181 \quad (2.10) \quad \lambda_\zeta(K) \stackrel{\text{def}}{=} \inf \left\{ \frac{- \int_{\mathcal{X}} f(x) G f(x) \pi(dx)}{\text{Var}_\pi(f) - \zeta}, \quad f \in \mathcal{A}, \text{Var}_\pi(f) > \zeta, \text{ and } \|f\|_* = 1 \right\}.$$

182 We have the analog of Lemma 2.1.

183 **Lemma 2.3.** *Suppose that the semigroup  $K$  satisfies (2.9). Let  $\nu(dx) = f(x)\pi(dx)$  be a  
 184 probability measure on  $\mathcal{X}$ , where  $f \in \mathcal{A}$ . Let  $\zeta \in [0, 1)$  be such that  $\lambda_\zeta(K) > 0$ . Then for all  
 185  $t \geq 0$  we have*

$$186 \quad \|\nu K_t - \pi\|_{\text{tv}}^2 \leq \text{Var}_\pi(K_t f) \leq \text{Var}_\pi(f) e^{-2\lambda_\zeta(K)t} + \zeta \|f\|_*^2.$$

187 *Proof.* See Section 5.3. ■

188 For  $\zeta = 0$ ,  $\lambda_\zeta(K)$  corresponds to the classical spectral gap of the semigroup and Lemma  
 189 2.3 is the classical exponential convergence of the semigroup. This result can be applied to  
 190 Langevin diffusion processes. Suppose that  $\mathcal{X} = \mathbb{R}^p$  equipped with the Lebesgue measure, and  
 191  $\pi(dx) = e^{-U(x)}/Z$ , for a function  $U : \mathbb{R}^p \rightarrow \mathbb{R}$  that is differentiable with Lipschitz gradient.  
 192 The Langevin diffusion process for  $\pi$  defines a reversible Markov semigroup with invariant dis-  
 193 tribution  $\pi$ . The convergence rate of the semigroup toward  $\pi$  is a key ingredient in the analysis  
 194 of several recent MCMC algorithms, including the unadjusted Langevin algorithm and sto-  
 195 chastic gradient Langevin dynamics ([Welling and Teh(2011), Raginsky et al.(2017)]). When  
 196  $U$  is convex, the semigroup is known to possess a spectral gap ([Bobkov(1999)]). Various exten-  
 197 sions beyond the convex case are also known and are well discussed in ([Bakry et al.(2008)]).  
 198 Lemma 2.3 offers another route, one that might be more effective when a good initial distri-  
 199 bution is available, and  $\pi$  has well-understood concentration properties. We leave the details  
 200 as possible future research.

201 **3. Application: mixing times of mixtures of Markov kernels.** To illustrate Lemma 2.1  
 202 we consider here the case where  $\mathcal{X} = \mathbb{R}^p$ , and  $\pi$  is a discrete mixture of log-concave densities  
 203 of the form

$$204 \quad (3.1) \quad \pi(dx) \propto \sum_{i \in \mathbb{I}} \pi(i, x) dx,$$

205 where  $\mathbb{I}$  is a nonempty finite set, and for  $i \in \mathbb{I}$ ,  $\pi(i, \cdot) : \mathbb{R}^p \rightarrow [0, \infty)$  is a measurable function.  
 206 As mentioned in the introduction, much is known on the computational complexity of various  
 207 MCMC algorithms for log-concave densities. However these results cannot be directly applied

208 to mixtures, since for instance a mixture of log-concave densities is not log-concave in general.  
 209 Sampling from mixtures is more challenging than sampling from log-concave densities. For  
 210 instance it is shown in [Ge et al.(2018)] that no polynomial-time MCMC algorithm exists to  
 211 sample from mixtures of densities with unequal covariance matrix, if the algorithm uses only  
 212 the marginal density of the mixture and its derivative. However this result does not cover the  
 213 most commonly used strategy to deal with mixtures, namely the Gibbs sampler.

214 Gibbs sampling type algorithms work with the joint distribution on  $I \times \mathcal{X}$  defined as

$$215 \quad (3.2) \quad \bar{\pi}(D \times B) = \frac{\sum_{i \in D} \int_B \pi(i, x) dx}{\sum_{i \in I} \int_{\mathcal{X}} \pi(i, x) dx}, \quad D \subseteq I, B \in \mathcal{B}.$$

216 Let  $\pi(i|x) \propto \pi(i, x)$  (resp.  $\pi(i) \propto \int_{\mathcal{X}} \pi(i, x) dx$ ) denote the implied conditional (resp.  
 217 marginal) distribution on  $I$ , and let  $\pi_i(dx) \propto \pi(i, x) dx$  be the implied conditional distribution  
 218 on  $\mathcal{X}$ . For each  $i \in I$ , let  $K_i$  be a transition kernel on  $\mathcal{X}$  with invariant distribution  $\pi_i$ . We  
 219 assume that  $K_i$  is reversible with respect to  $\pi_i$ , and ergodic (phi-irreducible and aperiodic).  
 220 We then consider the Markov kernel  $K$  defined as

$$221 \quad (3.3) \quad K(x, dy) \stackrel{\text{def}}{=} \sum_{i \in I} \pi(i|x) K_i(x, dy),$$

222 that is reversible with respect to  $\pi$  as in (3.1). In [Madras and Randall(2002)] the authors  
 223 developed a very nice lower bound on the spectral gap of  $K$  knowing the spectral gaps of the  
 224  $K_i$ 's. Their result goes as follows. Suppose that there exist  $\kappa > 0$ , and a graph on  $I$  such that  
 225 whenever there is an edge between  $i, j \in I$ , it holds

$$226 \quad (3.4) \quad \int_{\mathcal{X}} \min(\pi_i(x), \pi_j(x)) dx \geq \kappa.$$

227 If  $D(I)$  denotes the diameter of the graph thus defined<sup>1</sup>, Theorem 1.2 of [Madras and Randall(2002)]  
 228 says that

$$229 \quad (3.5) \quad \text{SpecGap}(K) \geq \frac{\kappa}{2D(I)} \min_{i \in I} \{\pi(i) \text{SpecGap}(K_i)\}.$$

230 The lower bound in (3.5) can be very small when  $I$  is large, particularly if some  $\pi(i)$   
 231 are exponentially small. We combine the approach in ([Madras and Randall(2002)]) with the  
 232 canonical path argument of ([Sinclair(1992), Diaconis and Stroock(1991)]) to develop a new  
 233 bound on the  $\zeta$ -spectral gap of  $K$ . We make the following assumption.

234 **H1.** *There exist  $I_0 \subseteq I$ , and  $\{B_i, i \in I_0\}$  a family of nonempty measurable subsets of  $\mathcal{X}$ ,*  
 235 *with the following property.*

- 236 1. *For each  $i \in I_0$ ,  $\pi_i(B_i) \geq 1/2$ .*
- 237 2. *There exist  $\kappa > 0$  and a connected graph  $\mathcal{G}$  on  $I_0$  such that*

$$238 \quad (3.6) \quad \int_{B_i \cap B_j} \min\left(\frac{\pi_i(x)}{\pi_i(B_i)}, \frac{\pi_j(x)}{\pi_j(B_j)}\right) dx \geq \kappa,$$

239 *whenever there is an edge in  $\mathcal{G}$  between  $i$  and  $j$ .*

---

<sup>1</sup>The diameter of a graph is the length (the number of edges) of the longest among all the shortest paths between all pairs of vertices.

240 One should view  $\cup_{i \in I_0} \{i\} \times \mathbf{B}_i$  as a subset of  $I \times \mathcal{X}$  that captures most of the probability mass  
 241 of  $\bar{\pi}$ . The graph  $\mathcal{G}$  captures the proximity between the conditional distributions. Indeed, (3.6)  
 242 implies that the total variation distance between the restriction of  $\pi_i$  to  $\mathbf{B}_i$  and the restriction  
 243 of  $\pi_j$  to  $\mathbf{B}_j$  is at most  $2(1 - \kappa)$ .

244 Since  $\mathcal{G}$  is assumed connected, for any distinct pair  $i, j \in I_0$  we can find and pick a path  
 245  $\gamma_{ij}$  that connects  $i$  and  $j$ . We call  $\gamma_{ij}$  the canonical path from  $i$  to  $j$ . The number of edges on  
 246  $\gamma_{ij}$  is denoted  $|\gamma_{ij}|$ . We then define

$$247 \quad (3.7) \quad \mathbf{m}_1 \stackrel{\text{def}}{=} \max_{\iota \in I_0} \sum_{i, j \in I_0: \gamma_{ij} \ni \iota} |\gamma_{ij}| \frac{\pi(i)\pi(j)}{\pi(\iota)},$$

248 where the summation is taken over all distinct pair  $(i, j)$  whose canonical path  $\gamma_{ij}$  goes through  
 249 node  $\iota$ . We define the local spectral gap of  $K_i$  as  $\text{SpecGap}_i(K_i) = \text{SpecGap}_{\mathbf{B}_i}(K_i)$ , where  
 250  $\text{SpecGap}_{\mathbf{B}_i}(K_i)$  is defined as in (2.6).

251 **Theorem 3.1.** *Let  $\pi$  as in (3.1), and  $K$  as in (3.3). Assume that H1 holds and  $K$  satisfies*  
 252 *(2.2) with some chosen pseudo-norm  $\|\cdot\|_*$ . Set  $\bar{\mathbf{B}} \stackrel{\text{def}}{=} \cup_{i \in I_0} \{i\} \times \mathbf{B}_i$  and assume that there*  
 253 *exists  $\zeta \in [0, 1)$  such that for any function  $f \in L^2_*(\pi)$  satisfying  $\|f\|_* = 1$ , it holds*

$$254 \quad (3.8) \quad 2 \int_{\bar{\mathbf{B}}} \int_{\bar{\mathbf{B}}^c} (f(y) - f(x))^2 \bar{\pi}(di, dx) \bar{\pi}(dj, dy) \\ 255 \quad \quad \quad + \int_{\bar{\mathbf{B}}^c} \int_{\bar{\mathbf{B}}^c} (f(y) - f(x))^2 \bar{\pi}(di, dx) \bar{\pi}(dj, dy) \leq \zeta,$$

256 where  $\bar{\mathbf{B}}^c \stackrel{\text{def}}{=} (I \times \mathcal{X}) \setminus \bar{\mathbf{B}}$ . Then

$$257 \quad (3.9) \quad \text{SpecGap}_\zeta(K) \geq \left( \frac{\kappa}{1 + 8\mathbf{m}_1} \right) \min_{i \in I_0} \text{SpecGap}_i(K_i).$$

260 *Proof.* See Section 5.4. ■

261 **Remark 3.2.** The condition (3.8) can be easily handled. For instance if  $\|\cdot\|_* = \|\cdot\|_{\pi, m}$  for  
 262 some  $m \in (2, \infty]$ , then by Holder's inequality the left hand side of (3.8) is easily bounded from  
 263 above by  $10\bar{\pi}(\bar{\mathbf{B}}^c)^{1-2/m}$ . In that case (3.8) holds if  $\bar{\mathbf{B}}$  satisfies  $\bar{\pi}(\bar{\mathbf{B}}) \geq 1 - (\zeta/10)^{1+2/(m-2)}$ .

264 Note that the constant  $\mathbf{m}_1$  satisfies

$$265 \quad (3.10) \quad \mathbf{m}_1 \leq \frac{D(I_0)}{\min_{i \in I_0} \pi(i)}.$$

266 Hence the bound in (3.9) improves on (3.5), even when  $\zeta = 0$ . In problems where an exact draw  
 267 from  $\pi(\cdot|x)$  is not available, the kernel  $K$  in (3.3) is not usable. In these cases it is typical to  
 268 replace those exact draws by MCMC. Theorem 3.1 can be extended to such settings. However  
 269 we will not pursue this here for lack of space.

270 **4. Example: analysis of a Gibbs sampler.** We consider the Bayesian treatment of a linear  
 271 regression problem with response variable  $z \in \mathbb{R}^n$ , and covariate matrix  $X \in \mathbb{R}^{n \times p}$ , with a  
 272 spike-and-slab prior distribution on the regression parameter  $\theta \in \mathbb{R}^p$  as in ([George and McCulloch(1997),  
 273 Narisetty and He(2014)]). More precisely, for some variable selection parameter  $\delta \in \Delta \stackrel{\text{def}}{=} \{0, 1\}^p$   
 274 and positive parameters  $\rho_0, \rho_1$ , we assume that the components of  $\theta$  are conditionally  
 275 independent, and  $\theta_j | \{\delta = 1\}$  has density  $\mathbf{N}(0, \rho_1^{-1})$ , and  $\theta_j | \{\delta = 0\}$  has density  $\mathbf{N}(0, \rho_0^{-1})$ ,  
 276 where  $\mathbf{N}(\mu, v^2)$  denotes the univariate Gaussian distribution with mean  $\mu$  and variance  $v^2$ .  
 277 We further assume that given  $\mathbf{q} \in (0, 1)$ , the prior distribution of  $\delta$  is a product of Bernoulli  
 278 with success probability  $\mathbf{q}$ , and restricted to be in  $\Delta_s \stackrel{\text{def}}{=} \{\delta \in \Delta : \|\delta\|_0 \leq s\}$ , for some sparsity  
 279 level  $s$  specified by the user. The resulting posterior distribution on  $\Delta \times \mathbb{R}^p$  is

$$280 \quad (4.1) \quad \Pi(\delta, d\theta|z) \propto \left( \frac{\mathbf{q}}{1-\mathbf{q}} \right)^{\|\delta\|_0} \mathbf{1}_{\Delta_s}(\delta) \frac{e^{-\frac{1}{2}\theta' D_{(\delta)}^{-1} \theta}}{\sqrt{\det(2\pi D_{(\delta)})}} e^{-\frac{1}{2\sigma^2} \|z - X\theta\|_2^2} d\theta,$$

281 where  $D_{(\delta)} \in \mathbb{R}^{p \times p}$  is a diagonal matrix with  $j$ -th diagonal element equal to  $\rho_1^{-1}$  if  $\delta_j = 1$ , and  
 282  $\rho_0^{-1}$  if  $\delta_j = 0$ . Note that we can always set  $s = p$ . The regression error  $\sigma$  is assumed known.  
 283 This model is very popular in the applications. Indeed, the posterior conditional distribution  
 284  $\Pi(\delta|\theta, z)$  is a product of independent Bernoulli distributions constrained to be  $s$ -sparse:

$$285 \quad (4.2) \quad \Pi(\delta|\theta, z) \propto \mathbf{1}_{\Delta_s}(\delta) \prod_{j=1}^p [\mathbf{q}_j]^{\delta_j} [1 - \mathbf{q}_j]^{1-\delta_j}, \quad \mathbf{q}_j \stackrel{\text{def}}{=} \frac{1}{1 + A e^{\frac{1}{2}(\rho_1 - \rho_0)\theta_j^2}}, \quad j = 1, \dots, p,$$

288 where  $A \stackrel{\text{def}}{=} (1 - \mathbf{q})\mathbf{q}^{-1} \sqrt{\rho_0/\rho_1}$ . We will assume that sampling from (4.2) is easy. This is the  
 289 case when  $s = p$  (by direct independent sampling), or when  $s$  is large (by a simple rejection  
 290 scheme). A Metropolis-Hastings scheme could also be used, but we will focus our analysis on  
 291 cases where an exact draw is made from (4.2). Given  $\delta$ , the conditional distribution of  $\theta$  given  
 292  $\delta$  is  $\mathbf{N}_p(m_\delta, \sigma^2 \Sigma_\delta)$ , with  $m_\delta$  and  $\Sigma_\delta$  given by

$$293 \quad (4.3) \quad m_\delta \stackrel{\text{def}}{=} \Sigma_\delta X' z \quad \text{and} \quad \Sigma_\delta \stackrel{\text{def}}{=} \left( X' X + \sigma^2 D_{(\delta)}^{-1} \right)^{-1}.$$

294 Put together these two conditional distributions yields a simple Gibbs sampling algorithm for  
 295 (4.1). We consider the following version that is modified so that the resulting Markov chain  
 296 is lazy as required by our theory.

---

[Algorithm 4] For some initial distribution  $\nu_0$  on  $\mathbb{R}^p$ , draw  $u_0 \sim \nu_0$ . Given  $u_0, \dots, u_k$  for some  $k \geq 0$ , draw independently  $I_{k+1} \sim \text{Ber}(0.5)$ .

1. If  $I_{k+1} = 0$ , set  $u_{k+1} = u_k$ .
  2. If  $I_{k+1} = 1$ ,
    - (a) Draw  $\delta \sim \Pi(\cdot|u_k, z)$  as given in (4.2), and
    - (b) draw  $u_{k+1} \sim \mathbf{N}_p(m_\delta, \sigma^2 \Sigma_\delta)$  as given in (4.3).
-

297 We analyze the mixing time of the marginal chain  $\{u_k, k \geq 0\}$  from Algorithm 4. As  
 298 easily seen,  $\{u_k, k \geq 0\}$  is a Markov chain with invariant distribution

$$299 \quad (4.4) \quad \Pi(d\theta|z) \propto \sum_{\delta \in \Delta_s} \binom{\mathbf{q}}{1-\mathbf{q}}^{\|\delta\|_0} \frac{e^{-\frac{1}{2}\theta' D_{(\delta)}^{-1} \theta}}{\sqrt{\det(2\pi D_{(\delta)})}} e^{-\frac{1}{2\sigma^2}\|z-X\theta\|_2^2} d\theta,$$

300 which is of the form (3.1), and with transition kernel

$$301 \quad (4.5) \quad K(u, d\theta) \stackrel{\text{def}}{=} \sum_{\omega \in \Delta} \Pi(\omega|u, z) \left[ \frac{1}{2} \delta_u(d\theta) + \frac{1}{2} \Pi(d\theta|\omega, z) \right],$$

302 which is of the form (3.3).

303 To proceed we introduce some notations. For  $\delta \in \Delta$ , and  $\theta \in \mathbb{R}^p$ , we write  $\theta_\delta$  as a short  
 304 for the component-wise product of  $\theta$  and  $\delta$ , and we define  $\delta^c \stackrel{\text{def}}{=} 1 - \delta$ , that is  $\delta_j^c = 1 - \delta_j$ ,  
 305  $1 \leq j \leq p$ . For a matrix  $A \in \mathbb{R}^{q \times p}$ ,  $A_\delta$  (resp.  $A_{\delta^c}$ ) denotes the matrix of  $\mathbb{R}^{q \times \|\delta\|_0}$  (resp.  
 306  $\mathbb{R}^{q \times (p - \|\delta\|_0)}$ ) obtained by keeping only the columns of  $A$  for which  $\delta_j = 1$  (resp.  $\delta_j = 0$ ).  
 307 When  $\delta = e_j$  (the  $j$ -th canonical unit vector of  $\mathbb{R}^p$ ) we write  $A_\delta$  (resp.  $A_{\delta^c}$ ) as  $A_j$  (resp.  $A_{-j}$ ).  
 308 For two elements  $\delta, \delta'$  of  $\Delta$ , we write  $\delta \supseteq \delta'$  to mean that  $\delta_j = 1$  whenever  $\delta'_j = 1$ . The support  
 309 of a vector  $u \in \mathbb{R}^p$  is the vector  $\text{supp}(u) \in \Delta$  such that  $\text{supp}(u)_j = 1$  if and only if  $|u_j| > 0$ .

310 An important role is played in the analysis by the matrices

$$311 \quad L_\delta \stackrel{\text{def}}{=} I_n + \frac{1}{\sigma^2} X D_{(\delta)} X',$$

312 and the coherence of  $X$  defined as

$$313 \quad \mathcal{C}(s) \stackrel{\text{def}}{=} \max_{\delta \in \Delta_s} \max_{j \neq \ell} \frac{|X'_j L_\delta^{-1} X_\ell|}{\sqrt{n \log(p)}}.$$

314 We will make the assumption that  $\mathcal{C}(s)$  does not grow with  $p$ . It can be easily checked that  
 315 if the columns of  $X$  are orthogonal then  $\mathcal{C}(s) = 0$ . Furthermore, it can be shown that if  $X$   
 316 is a realization of random matrix with i.i.d. standard Gaussian entries, then and provided that  
 317  $n \geq As^2 \log(p)$ , it holds  $\mathcal{C}(s) \leq c$  for some absolute constants  $c, A$ . We refer the reader to the  
 318 Appendix for details. We make the following regularity assumption on the matrix  $X$ .

319 **H2.** 1. *The matrix  $X$  is non-random and normalized such that*

$$320 \quad (4.6) \quad \|X_j\|_2^2 = n, \quad j = 1, \dots, p.$$

321 *Furthermore, there exists an integer  $s_0 \in \{1, \dots, p-1\}$ , such that*

$$322 \quad \lambda \stackrel{\text{def}}{=} \min_{\delta: \|\delta\|_0 \leq s_0} \inf \left\{ \frac{v' (X'_{\delta^c} L_\delta^{-1} X_{\delta^c}) v}{n \|v\|_2^2}, v \in \mathbb{R}^{p-\|\delta\|_0}, 0 < \|v\|_0 \leq s_0 \right\} > 0.$$

323 *Remark 4.1.* The matrix  $L_\delta^{-1}$  can be loosely interpreted as the projector on the orthogonal  
 324 of the space spanned by the columns of  $X_\delta$ . Therefore, H2 rules out settings where a small  
 325 number of columns of  $X$  have the same column span as the column span of  $X$ . Indeed signal  
 326 recovery becomes nearly impossible in such settings. It can be shown that if  $X$  is a random  
 327 matrix with i.i.d. standard Gaussian entries then  $\lambda > 0$  for  $s_0$  of order  $n/\log(p)$ . We refer the  
 328 reader to the Appendix for details.  $\square$

330 We also make some very mild assumptions pertaining to the prior parameters and to the  
 331 existence of a true model.

332 **H3.** 1. *There exists a true value of the parameter  $\theta_\star \in \mathbb{R}^p$  with sparsity support*  
 333  *$\delta_\star \in \Delta_s$ , with  $\|\delta_\star\|_0 = s_\star$ , such that  $p^{s_\star} \Pi(\delta_\star|z) \geq 1$ .*  
 334 2. *For some constant  $u > 0$ , the prior parameter  $\mathbf{q}$  satisfies*

$$335 \quad (4.7) \quad \frac{\mathbf{q}}{1 - \mathbf{q}} = \frac{1}{p^u}.$$

336 3. *The prior parameters  $\rho_0, \rho_1$  satisfy*

$$337 \quad (4.8) \quad 0 < \rho_1 < \rho_0, \quad \sigma^2 \rho_1 \leq \left(1 - \frac{\rho_1}{\rho_0}\right) n, \quad \text{and} \quad \sqrt{1 + \frac{ns}{\sigma^2 \rho_1}} \leq p^a,$$

338 *for some absolute constant  $a > 0$ .*

339 The last two parts of Condition (4.8) are easily satisfied and are imposed mostly to obtain  
 340 simple mathematical formulas. For some constant  $c_0 > 0$ , we introduce the event

$$341 \quad \mathcal{E}_0 \stackrel{\text{def}}{=} \left\{ z \in \mathbb{R}^n : \max_{\delta \in \Delta_s} \sup_{1 \leq j \leq p} \frac{1}{\sigma} |\langle L_\delta^{-1} X_j, z - X \theta_\star \rangle| \leq \sqrt{c_0 n \log(p)} \right\},$$

342 We note if  $z \sim \mathbf{N}(X \theta_\star, \sigma^2 I_n)$ , and  $\|X_j\|_2 \leq \sqrt{n}$ , then the event  $z \in \mathcal{E}_0$  holds with high  
 343 probability, with  $c_0 = 2(s+1)$ .

344 **Theorem 4.2.** *Suppose that H2-H3 hold. Fix  $\zeta_0 \in (0, 1)$ . Suppose that  $s$ , the sparsity level  
 345 of the posterior distribution (4.1) is chosen such that  $0 < s \leq s_0$  with  $s_0$  as in H2, and  
 346 Algorithm 4 is initialized from  $\nu_0 = \Pi(\cdot|\delta^{(i)}, z)$ , for some arbitrary  $\delta^{(i)} \in \Delta_s$ . Take  $z \in \mathcal{E}_0$ ,  
 347 suppose that we choose  $u$  large enough such that*

$$348 \quad (4.9) \quad u > 2 \max \left( 2, \frac{\varrho}{\lambda} \right), \quad \text{where} \quad \varrho \stackrel{\text{def}}{=} (\sigma \sqrt{c_0} + \|\theta_\star\|_1 \mathcal{C}(s))^2,$$

349 *and the sample size  $n$  satisfies*

$$350 \quad (4.10) \quad n \geq \frac{A_0 u \sigma^2 s_\star \log(p)}{\lambda^2 \underline{\theta}_\star^2}, \quad \text{where} \quad \underline{\theta}_\star \stackrel{\text{def}}{=} \min_{j: \delta_\star j = 1} |\theta_{\star j}|,$$

351 *for some absolute constant  $A_0$ . Set*

$$352 \quad \lambda_1 \stackrel{\text{def}}{=} \min_{1 \leq j \leq p} \min_{\delta \in \Delta_s} \frac{X_j' L_\delta^{-1} X_j}{n}.$$

353 Then there exists a constant  $A_1$  that does not depend on  $n, p$  nor  $\zeta_0$  such that for all

354

$$355 \quad (4.11) \quad N \geq A_1 s \left[ \log \left( \frac{1}{\zeta_0} \right) + \frac{su(1 + \|\theta_\star\|_\infty^2)n}{\sigma^2 \lambda} \right] \times \max \left( 1, \sqrt{\frac{n}{\sigma^2 \rho_0}} \right) \\ 356 \quad \quad \quad \times \max \left( 1, e^{\frac{1}{4\sigma^2}(8\sigma^2 \rho_0 - n\lambda_1)} \right) \times p^{\frac{\rho_0}{n} \frac{2g}{\lambda_1^2}}, \\ 357$$

358 we have

359

$$\|\nu_0 K^N - \Pi(\cdot|z)\|_{\text{tv}} \leq \zeta_0.$$

360 *Proof.* See Section 5.5. ■

361 We note that our condition (4.9) is analogous to Condition C of [Yang et al.(2016)]. The  
362 main term in the bound (4.11) is

$$363 \quad \max \left( 1, e^{\frac{1}{4\sigma^2}(8\sigma^2 \rho_0 - n\lambda_1)} \right) p^{\frac{\rho_0}{n} \frac{2g}{\lambda_1^2}},$$

364 which highlights the important impact of the prior parameter  $\rho_0$  on the mixing of the algo-  
365 rithm. If  $\rho_0$  is chosen as  $\rho_0 \leq n\lambda_1/(8\sigma^2)$ , then by (4.11), the mixing time scales as  $O(p^{\rho/\lambda_1})$ .  
366 Note that the ratio  $\rho/\lambda_1$  depends mainly on the correlation between the columns of  $X$ . Our  
367 simulation results indeed confirm that dependence of the mixing time on  $X$ , however the  
368 polynomial scaling  $O(p^{\rho/\lambda_1})$  predicted by the theorem may be conservative.

369 In contrast, if  $\rho_0 > n\lambda_1/(8\sigma^2)$  the bound predicts a mixing time that scales as  $O(e^{2\rho_0} p^{\frac{\rho_0}{n} \frac{2g}{\lambda_1^2}})$ ,  
370 which is worst than  $O(e^n p^{\rho/\lambda_1})$ . This said, it is important to add that (4.11) is an upper bound  
371 on the mixing time which may not be tight, and as such does not prove slow mixing.

372 We contrast these findings with the posterior contraction properties of the posterior dis-  
373 tribution. According to [Narisetty and He(2014)], as  $n, p \rightarrow \infty$ , we need to let  $\rho_0$  grow faster  
374 than  $n$ , and let  $\rho_1$  be of order  $n/p^2$  in order to guarantee posterior contraction of  $\Pi$ . And  
375 in their simulation section these authors suggest using  $\rho_0 = 10n/\sigma^2$  (although it is unclear  
376 whether posterior contraction holds in that regime). In these regimes our results suggest that  
377 the mixing time of Algorithm 4 grows faster than  $O(e^n p^{\rho/\lambda_1})$ . This description matches well  
378 with our numerical experiments. But again (4.11) is only an upper bound on the mixing time,  
379 and as such does not establish slow mixing.

380 Note that when posterior contraction holds the posterior distribution assigns increasingly  
381 small probability to  $\{\delta : \delta \not\supseteq \delta_\star\}$ . Hence a chain that starts in  $\{\delta : \delta \supseteq \delta_\star\}$  may have  
382 markedly different mixing time than what is predicted by Theorem 4.2. To formalize this,  
383 we shall focus on the unconstrained case where  $s = p$  in (4.1). We formalize the posterior  
384 contraction as follows. Given  $k \geq 0$ , we define

$$385 \quad \mathcal{D}_k \stackrel{\text{def}}{=} \{\delta \in \Delta : \delta \supseteq \delta_\star, \|\delta\|_0 \leq \|\delta_\star\|_0 + k\},$$

386 which collects models that contain the true model  $\delta_\star$  and have at most  $k$  false-positives, and

387 we introduce the event  
388

$$389 \quad \mathcal{E} \stackrel{\text{def}}{=} \left\{ z \in \mathbb{R}^n : \Pi(\mathcal{D}_k|z) \geq 1 - \frac{1}{p^{\frac{u}{2}(k+1)}}, \text{ for all } k \geq 0, \right.$$

$$390 \quad \left. \text{and } \max_{\delta \supseteq \delta_\star: \|\delta\|_0 \leq s_0} \sup_{1 \leq j \leq p} \frac{1}{\sigma} |\langle L_\delta^{-1} X_j, z - X\theta_\star \rangle| \leq \sqrt{c_0 n \log(p)} \right\},$$

391

392 for some constant  $c_0$ . We will say that posterior contraction holds when  $z \in \mathcal{E}$ . We will not  
393 directly establish this property. However several existing works suggest that this description  
394 of the posterior contraction of  $\Pi(\cdot|z)$  holds. For instance under similar assumptions as above,  
395 [Narisetty and He(2014)] show that  $\Pi(\mathcal{D}_0|Z) \geq 1 - \frac{a_1}{p^{a_2}}$  with high-probability for positive con-  
396 stants  $a_1, a_2$ . And [Atchade and Bhattacharyya(2018)] shows that  $z \in \mathcal{E}$  with high probabiity  
397 for a slightly modified version of the posterior distribution (4.1).

398 **Theorem 4.3.** *Assume H2-H3 and  $s = p$  in (4.1). Fix  $\zeta_0 \in (0, 1)$ . Suppose that Algorithm  
399 4 is initialized from  $\nu_0 = \Pi(\cdot|\delta^{(i)}, z)$ , for some  $\delta^{(i)} \in \mathcal{D}_{(s_0-s_\star)}$  such that  $FP \stackrel{\text{def}}{=} \|\delta^{(i)}\|_0 - s_\star$   
400 satisfies*

$$401 \quad (4.12) \quad FP \leq \frac{u}{4(u+a)}(k+1) + \frac{\log\left(\frac{80}{\zeta_0^2}\right)}{2(u+a)\log(p)},$$

402 for some integer  $k \leq s_0 - s_\star$ . Suppose also that (4.9) and (4.10) hold. Then there exists a  
403 constant  $A$  that does not depend on  $n, p$  nor  $\zeta_0$  such that for all  $z \in \mathcal{E}$ , and all

$$404 \quad (4.13) \quad N \geq A FP [\log(\zeta_0^{-1}) + FP u \log(p)] p^{\frac{2\rho_0}{n} \frac{\rho}{\lambda_1^2}},$$

405 we have

$$406 \quad \|\nu_0 K^N - \Pi(\cdot|z)\|_{\text{tv}} \leq \zeta_0.$$

407 *Proof.* See Section 5.6. ■

408 Condition (4.12) restricts the number of false-positives of the initial model  $\delta^{(i)}$  compared  
409 to  $s_0$ . This condition can be relaxed if the contraction of  $\pi$  on  $\mathcal{D}_k$  is faster than the polynomial  
410 form assumed in the event  $\mathcal{E}$ .

411 Theorem 4.3 suggests that when posterior contraction holds ( $z \in \mathcal{E}$ ), the mixing time  
412 of Algorithm 4 with a good initialization is less sensitive to large values of  $\rho_0$  (the term  
413  $e^{\frac{1}{4\sigma^2}(8\sigma^2\rho_0 - n\lambda_1)}$  no longer appear in (4.13)). For instance with  $\rho_0 = n\lambda_1/2$  the mixing time is  
414 at most  $O(FP^2 p^{\rho/\lambda_1})$ , which is better  $O(e^n p^{\rho/\lambda_1})$ .

415 One clear roadblock toward the practical use of this result is finding the initial  $\delta^{(i)}$  such  
416 that  $\delta^{(i)} \supseteq \delta_\star$ . In practice various frequentist estimators such as the lasso can be used.  
417 At least in a high signal-to-noise-ratio setting the lasso estimator is known to contain the  
418 true model under mild assumptions (similar to H2). We refer the reader for instance to  
419 ([Meinshausen and Yu(2009)]).

420 One of the first paper that analyzes the mixing times of MCMC algorithm in high-  
421 dimensional linear regression models and highlights fast/slow mixing behaviors is [Yang et al.(2016)]. ■

422 Their posterior distribution is slightly different from what we looked at in this work. Specifi-  
 423 cally [Yang et al.(2016)] applied a Metropolized-Gibbs sampler to the marginal distribution of  
 424  $\delta$ , whereas we consider here a Gibbs sampler applied to the joint distribution of  $(\delta, \theta)$ . These  
 425 authors show that in general their sampler has a mixing time that is exponential in  $p$  unless  
 426 the state space is restricted to models  $\delta$  for which  $\|\delta\|_0 \leq s$  for some threshold  $s$ , in which case  
 427 the worst-case mixing time is  $O(s^2 np \log(p))$ . To the extent that our bound in Theorem 4.2 is  
 428 tight, the better rate obtained by these authors can perhaps be interpreted as the positive  
 429 effect of marginalization and collapsing in Gibbs sampling ([Liu(1994)]).

430 **4.1. Numerical illustrations.** We illustrate some of the conclusions with the following sim-  
 431 ulation study. We consider a linear regression model with Gaussian noise  $\mathbf{N}(0, \sigma^2)$ , where  $\sigma^2$  is  
 432 set to 1. We experiment with sample size  $n = p$ , and dimension  $p \in \{500, 1000, 2000, 3000, 4000\}$ .  
 433 We take  $X \in \mathbb{R}^{n \times p}$  as a random matrix with i.i.d. rows drawn from  $\mathbf{N}_p(0, \Sigma)$  under two sce-  
 434 narios. A low coherence setting where  $\Sigma = I_p$ , and a high coherence where  $\Sigma_{ij} = 0.9^{|j-i|}$ .  
 435 After sampling, we normalized the columns of  $X$  to each have norm  $\sqrt{n}$ . We fix the number  
 436 of non-zero coefficients to  $s_\star = 10$ , and  $\delta_\star$  is given by

$$437 \quad \delta_\star = (\underbrace{1, \dots, 1}_{10}, \underbrace{0, \dots, 0}_{p-10}).$$

438 The non-zero coefficients of  $\theta_\star$  are uniformly drawn from  $(-a - 1, -a) \cup (a, a + 1)$ , where

$$439 \quad a = 4\sqrt{\frac{\log(p)}{n}}.$$

440 We use the following prior parameters values

$$441 \quad u = 2, \quad \rho_1 = \frac{n}{p^{2.1}}, \quad \rho_0 \in \left\{ \frac{n}{\sigma^2}, \frac{n^{1.5}}{\sigma^2} \right\}.$$

442 These scalings of  $\rho_0$  and  $\rho_1$  roughly matches the recommendations of [Narisetty and He(2014)]  
 443 to get posterior contraction of  $\Pi(\cdot|z)$ . We use an initial distribution  $\nu_0 = \Pi(\cdot|\delta^{(i)}, z)$ , where  
 444  $\delta^{(i)}$  is such that  $\|\delta^{(i)} - \delta_\star\|_0 = 2p/10$ , with two scenarios. A scenario FN (false negative),  
 445 where 5 out of 10 of the true positive of  $\delta_\star$  are set to 0, and a scenario no FN, where  $\delta^{(i)}$  has  
 446 only false-positives. To monitor the mixing, we compute the sensitivity and the precision at  
 447 iteration  $k$  as  
 448

$$449 \quad \text{SEN}_k = \frac{1}{s_\star} \sum_{j=1}^p \mathbf{1}_{\{|\delta_{k,j}|>0\}} \mathbf{1}_{\{|\delta_{\star,j}|>0\}}, \quad \text{PREC}_k = \frac{\sum_{j=1}^p \mathbf{1}_{\{|\delta_{k,j}|>0\}} \mathbf{1}_{\{|\delta_{\star,j}|>0\}}}{\sum_{j=1}^p \mathbf{1}_{\{|\delta_{k,j}|>0\}}}.$$

451 We empirically measure the mixing time of the algorithm as the first time  $k$  where both  $\text{SEN}_k$   
 452 and  $\text{PREC}_k$  reach 1, truncated to  $2 \times 10^4$  – that is we stop any run that has not mixed by  
 453 20000 iterations. For the sampler of [Yang et al.(2016)], we stop any run that has not mixed  
 454 by  $10^5$  iterations. The average empirical mixing time thus obtained (based on 50 independent  
 455 MCMC replications) are presented in Table 1 and Table 2.

456 We can make the following observations.

		$p = 500$	$p = 1000$	$p = 2000$	$p = 3000$	$p = 4000$
FN	$\rho_0 = n$	866.3(3, 204)	423.6(2, 735)	147.1(575)	> 437.3	> 871.0
	$\rho_0 = n^{1.5}$	> 11, 125.8	> 13, 662.6	> 13, 2371.6	> 15, 948.0	> 16237.3
	Yang et al.	5, 244.2(1, 379)	12, 208.5(2, 463)	27, 617.6(5, 803)	43, 821.9(6, 453)	54, 697.9(5, 611)
no FN	$\rho_0 = n$	1(0)	1(0)	1(0)	1(0)	1(0)
	$\rho_0 = n^{1.5}$	30.9(81)	43.7(55)	123.2(251)	241.2(535)	215.3(250)
	Yang et al.	5, 191.0(1, 503)	11, 975.9(2, 769)	26, 877.8(4, 786)	42, 285.7(8, 721)	56, 264.3(10, 362)

**Table 1**

Average empirical mixing time of the samplers in a low-coherence setting. Based on 50 simulation replications. The numbers in parenthesis are standard errors. The notation  $> a$  means that some (or all) of the replicated mixing times have been truncated.

		$p = 500$	$p = 1000$	$p = 2000$	$p = 3000$	$p = 4000$
FN	$\rho_0 = n$	> 20, 000	> 19, 200	> 18, 400	> 17, 870	> 19129.1
	$\rho_0 = n^{1.5}$	> 20, 000	> 20, 000	> 20, 000	> 20, 000	> 20, 000
	Yang et al.	> 100, 000	> 91, 177	> 75, 373	> 83, 246	> 84, 972
no FN	$\rho_0 = n$	> 880.1	> 1, 200.1	> 400.9	> 800.96	> 900.1
	$\rho_0 = n^{1.5}$	> 416.8	> 1, 246.2	> 874.2	> 425.2	> 313.6
	Yang et al.	> 98, 067	> 87, 424	> 73, 253	> 77, 902	> 82, 205

**Table 2**

Average empirical mixing time of the samplers in a high-coherence setting. Based on 50 simulation replications. The numbers in parenthesis are standard errors. The notation  $> a$  means that some (or all) of the replicated mixing times have been truncated.

- 457 1. There is sharp difference in behavior between the low and high coherence settings.  
458 2. As predicted by our theory, Algorithm 4 mixes better when there is no false-negative  
459 in the initialization. The algorithm of [Yang et al.(2016)] seems impervious to the  
460 initialization. It should be noted in comparing the two algorithms, that an iteration  
461 of the algorithm of [Yang et al.(2016)] costs roughly  $p$  times less than an iteration of  
462 Algorithm 4.  
463 3. The third observation that can be drawn from the results is that when there are false-  
464 negatives, Algorithm 4 mixes better with  $\rho_0 = n/\sigma^2$ , compared to  $\rho_0 > n/\sigma^2$ , as  
465 predicted by our result. The difference is less noticeable in the high-coherence setting.  
466 This observation is also explained by our bound, since in a high-coherence setting, the  
467 parameter  $\rho$  is expected to be large. Another observation here is that when there are  
468 false-negatives in the initialization, the mixing time becomes highly variable (several  
469 runs have hit the wallclock).  
470 4. Finally, we notice that the theory of [Yang et al.(2016)] does not fully describe the  
471 behavior of their algorithm, as we see a significant degradation of performance in their  
472 algorithm with high coherence design matrices, which cannot be clearly explained by  
473 their result.

474 Overall, based on our theoretical analysis and the simulation study, our recommendation  
475 when using Algorithm 4 is to set  $\rho_0 = n/\sigma^2$ , and to the extent possible to use the lasso sparsity  
476 structure as initialization (or some other similar high-dimensional frequentist estimator).

477 **5. Proofs.** The proof of Theorem 3.1 relies on the following lemma due to [Madras and Randall(2002)]. ■  
 478 For a proof see their inequality (47). A direct argument by coupling can also be easily con-  
 479 structed.

480 **Lemma 5.1.** *Let  $\nu(dx) = f_\nu(x)dx$ ,  $\mu(dx) = f_\mu(x)dx$  be two probability measures on some*  
 481 *measurable space with reference measure  $dx$ , such that  $\int \min(f_\mu(x), f_\nu(x))dx > \epsilon$  for some*  
 482  *$\epsilon > 0$ . Then for any measurable function  $h$  such that  $\int h^2(x)\nu(dx) < \infty$  and  $\int h^2(x)\mu(dx) <$*   
 483  *$\infty$ , we have*

$$484 \int (h(y) - h(x))^2 \mu(dy) \nu(dx)$$

$$485 \leq \frac{2 - \epsilon}{2\epsilon} \left[ \int (h(y) - h(x))^2 \mu(dy) \mu(dx) + \int (h(y) - h(x))^2 \nu(dy) \nu(dx) \right].$$

488 **5.1. Proof of Lemma 2.1.** We first note that if a probability measure  $\nu$  is absolutely  
 489 continuous with respect to  $\pi$  with Radon-Nikodym derivative  $f_\nu$ , then for any  $A \in \mathcal{B}$ ,

$$490 \nu K(A) = \int \nu(dx) K(x, A) = \int \int f_\nu(x) \mathbf{1}_A(y) \pi(dx) K(x, dy)$$

$$491 = \int \int \mathbf{1}_A(x) f_\nu(y) \pi(dx) K(x, dy) = \int_A \pi(dx) \int K(x, dy) f_\nu(y),$$

492 where the third equality uses the reversibility of  $K$ . This calculation says that  $\nu K$  is also  
 493 absolutely continuous with respect to  $\pi$  with Radon-Nikodym derivative  $x \mapsto K f_\nu(x) \stackrel{\text{def}}{=} \int K(x, dy) f_\nu(y)$ . More generally  $\frac{d(\nu K^n)}{d\pi}(\cdot) = K^n f_\nu(\cdot)$ , and

$$495 \|\nu K^n - \pi\|_{\text{tv}}^2 = \left( \int \left| \frac{d(\nu K^n)}{d\pi}(x) - 1 \right| \pi(dx) \right)^2$$

$$496 = \left( \int |K^n f_\nu(x) - 1| \pi(dx) \right)^2$$

$$497 \leq \|K^n f_\nu - 1\|_{2, \pi}^2$$

$$498 (5.1) \quad = \text{Var}_\pi(K^n f_\nu).$$

499 Take  $f \in L^2(\pi)$ . Since  $\pi(f) = \pi(Kf)$ , we have

$$500 (5.2) \quad \text{Var}_\pi(Kf) - \text{Var}_\pi(f) = \langle Kf, Kf \rangle_\pi - \langle f, f \rangle_\pi = -\frac{1}{2} \int \int (f(y) - f(x))^2 \pi(dx) K^2(x, dy),$$

502

503 where the last equality exploits the reversibility of  $K$ . By the lazyness of  $K$  we have

$$504 \int \int (f(y) - f(x))^2 \pi(dx) K^2(x, dy) \geq \int \int (f(y) - f(x))^2 \pi(dx) K(x, dy).$$

505 A proof of this statement is given for instance in [Montenegro and Tetali(2006)] (Equation  
 506 2.12). Using the last display together with (5.2), and the definition of  $\mathcal{E}(f, f)$ , we conclude  
 507 that for all  $f \in L^2(\pi)$ ,

$$508 \quad (5.3) \quad \text{Var}_\pi(Kf) \leq \text{Var}_\pi(f) - \mathcal{E}(f, f).$$

509 Fix  $\zeta \in (0, 1)$ , and take  $f \in L^2_\star(\pi)$ . Suppose that  $\|f\|_\star > 0$ . If  $\text{Var}_\pi(f) \leq \zeta \|f\|_\star^2$ , then, by  
 510 (5.3),  $\text{Var}_\pi(Kf) \leq \min(\text{Var}_\pi(f), \zeta \|f\|_\star^2)$ . But if  $\text{Var}_\pi(f) > \zeta \|f\|_\star^2$ , then by (5.3),

$$\begin{aligned} 511 \quad \text{Var}_\pi(Kf) &\leq \text{Var}_\pi(f) - \|f\|_\star^2 \mathcal{E}\left(\frac{f}{\|f\|_\star}, \frac{f}{\|f\|_\star}\right) \\ 512 &\leq \text{Var}_\pi(f) - \|f\|_\star^2 \text{SpecGap}_\zeta(K) \left(\text{Var}_\pi\left(\frac{f}{\|f\|_\star}\right) - \frac{\zeta}{2}\right), \\ 513 &\leq \text{Var}_\pi(f) (1 - \text{SpecGap}_\zeta(K)) + \frac{\zeta}{2} \|f\|_\star^2 \text{SpecGap}_\zeta(K). \end{aligned}$$

514 Note also that if  $\|f\|_\star = 0$ , then  $\text{Var}_\pi(f) = 0$  by the listed properties of  $\|\cdot\|_\star$ , and  $\text{Var}_\pi(Kf) = 0$   
 515 by (5.3), so that the last display continue to hold. We conclude that for all  $f \in L^2_\star(\pi)$ ,

$$516 \quad \text{Var}_\pi(Kf) \leq \text{Var}_\pi(f) (1 - \text{SpecGap}_\zeta(K)) + \zeta \|f\|_\star^2 \text{SpecGap}_\zeta(K).$$

517 Given that  $Kf \in L^2_\star(\pi)$  for all  $f \in L^2_\star(\pi)$ , we can iterate the above inequality to deduce that  
 518 for all  $f \in L^2_\star(\pi)$ , and for all  $n \geq 1$ ,

$$\begin{aligned} 519 \quad \text{Var}_\pi(K^n f) &\leq \text{Var}_\pi(f) (1 - \text{SpecGap}_\zeta(K))^n \\ 520 &\quad + \zeta \text{SpecGap}_\zeta(K) \sum_{j \geq 0} (1 - \text{SpecGap}_\zeta(K))^j \|K^{n-j-1} f\|_\star^2 \\ 521 &\quad \leq \text{Var}_\pi(f) (1 - \text{SpecGap}_\zeta(K))^n + \zeta \|f\|_\star^2. \end{aligned}$$

523

524 Now, if  $\pi_0 = f_0 \pi$ , the last display combined with (5.1) implies that

$$525 \quad \|\pi_0 K^n - \pi\|_{\text{tv}}^2 \leq \text{Var}_\pi(K^n f_0) \leq \text{Var}_\pi(f_0) (1 - \text{SpecGap}_\zeta(K))^n + \zeta \|f_0\|_\star^2,$$

526 as claimed. □

528 **5.2. Proof Lemma 2.2.** Take  $f : \mathcal{X} \rightarrow \mathbb{R}$  such that  $\text{Var}_\pi(f) > \zeta$ , and  $\|f\|_\star = \|f\|_{m, \pi} = 1$ .  
 529 We have

530

$$\begin{aligned} 531 \quad 2\text{Var}_\pi(f) &= \int_{\mathcal{X}_0} \int_{\mathcal{X}_0} (f(y) - f(x))^2 \pi(dx) \pi(dy) \\ 532 &\quad + 2 \int_{\mathcal{X}_0} \int_{\mathcal{X} \setminus \mathcal{X}_0} (f(y) - f(x))^2 \pi(dx) \pi(dy) + \int_{\mathcal{X} \setminus \mathcal{X}_0} \int_{\mathcal{X} \setminus \mathcal{X}_0} (f(y) - f(x))^2 \pi(dx) \pi(dy). \end{aligned}$$

533

534 Using the convexity inequality  $(a + b)^2 \leq 2a^2 + 2b^2$ , and Holder's inequality,

535

$$\begin{aligned}
 536 \quad & \int_{\mathcal{X}_0} \int_{\mathcal{X} \setminus \mathcal{X}_0} (f(y) - f(x))^2 \pi(dx) \pi(dy) \\
 537 \quad & \leq 2\pi(\mathcal{X}_0) \int_{\mathcal{X} \setminus \mathcal{X}_0} f(x)^2 \pi(dx) + 2\pi(\mathcal{X} \setminus \mathcal{X}_0) \int_{\mathcal{X}_0} f(x)^2 \pi(dx) \\
 538 \quad & \leq 2\pi(\mathcal{X}_0) \pi(\mathcal{X} \setminus \mathcal{X}_0)^{1 - \frac{2}{m}} \|f\|_{m,\pi}^2 + 2\pi(\mathcal{X} \setminus \mathcal{X}_0) \|f\|_{m,\pi}^2 \leq 4\pi(\mathcal{X} \setminus \mathcal{X}_0)^{1 - \frac{2}{m}}.
 \end{aligned}$$

539

540 With similar calculation,

$$541 \quad \int_{\mathcal{X} \setminus \mathcal{X}_0} \int_{\mathcal{X} \setminus \mathcal{X}_0} (f(y) - f(x))^2 \pi(dx) \pi(dy) \leq 4\pi(\mathcal{X} \setminus \mathcal{X}_0) \pi(\mathcal{X} \setminus \mathcal{X}_0)^{1 - \frac{2}{m}} \leq 2\pi(\mathcal{X} \setminus \mathcal{X}_0)^{1 - \frac{2}{m}}.$$

542 Using  $\pi(\mathcal{X}_0) \geq (\zeta/10)^{1+2/(m-2)}$ , we get

$$543 \quad 2(\text{Var}_\pi(f) - \frac{\zeta}{2}) \geq \int_{\mathcal{X}_0} \int_{\mathcal{X}_0} \pi(dx) \pi(dy) (f(y) - f(x))^2.$$

544 Hence

$$545 \quad \frac{\mathcal{E}(f, f)}{\text{Var}_\pi(f) - \frac{\zeta}{2}} \geq \frac{\int_{\mathcal{X}_0} \int_{\mathcal{X}_0} \pi(dx) K(x, dy) (f(y) - f(x))^2}{\int_{\mathcal{X}_0} \int_{\mathcal{X}_0} \pi(dx) \pi(dy) (f(y) - f(x))^2} \geq \text{SpecGap}_{\mathcal{X}_0}.$$

546 The statement bound easily follows.  $\square$

548 **5.3. Proof of Lemma 2.3.** Take  $f \in \mathcal{A}$ . Without any loss of generality we assume that  
 549  $\pi(f) = 0$ . Then

$$550 \quad (5.4) \quad \frac{d}{dt} \text{Var}_\pi(K_t f) = \frac{d}{dt} \int_{\mathcal{X}} (K_t f)^2(x) \pi(dx) = 2 \int_{\mathcal{X}} K_t f(x) G K_t f(x) \pi(dx).$$

551 Suppose that  $\|K_t f\|_\star > 0$ . If  $\text{Var}_\pi(K_t f / \|K_t f\|_\star) > \zeta$ , then from (5.4) and the definition of  
 552  $\lambda_\zeta(K)$ ,

553

$$\begin{aligned}
 554 \quad (5.5) \quad \frac{d}{dt} \text{Var}_\pi(K_t f) & \leq -2\|K_t f\|_\star^2 \lambda_\zeta(K) \left( \text{Var}_\pi \left( \frac{K_t f}{\|K_t f\|_\star} \right) - \zeta \right) \\
 555 \quad & \leq -2\lambda_\zeta(K) \text{Var}_\pi(K_t f) + 2\zeta \lambda_\zeta(K) \|K_t f\|_\star^2.
 \end{aligned}$$

557 However, if  $\text{Var}_\pi(K_t f / \|K_t f\|_\star) \leq \zeta$ , we see that the right-hand side of (5.5) is nonnegative,  
 558 whereas from (5.4) and the properties of the generator we see that the left-hand side of (5.5)  
 559 is nonpositive. Note also that (5.5) continue to hold when  $\|K_t f\|_\star = 0$ . Hence for all  $f \in \mathcal{A}$ ,  
 560 and for all  $t \geq 0$ , we have

$$561 \quad (5.6) \quad \frac{d}{dt} \text{Var}_\pi(K_t f) \leq -2\lambda_\zeta(K) \text{Var}_\pi(K_t f) + 2\zeta \lambda_\zeta(K) \|f\|_\star^2.$$

562 The lemma then follows from Gronwall's lemma. More precisely, set  $\alpha = \zeta \|f\|_*^2$ ,  $\beta = 2\lambda_\zeta(K)$ ,  
 563 and  $u(t) = \text{Var}_\pi(K_t f)$ . Hence (5.6) reads  $u'(t) \leq -\beta u(t) + \alpha\beta$ . Setting  $v(t) = e^{-\beta t}$ , we have

$$564 \quad \frac{d}{dt} \left( \frac{u(t)}{v(t)} \right) = \frac{u'(t)v(t) - v'(t)u(t)}{v(t)^2} = \frac{u'(t) + \beta u(t)}{v(t)} \leq \alpha\beta e^{\beta t}.$$

565 Integrating both sides yields the stated bound.  $\square$

567 **5.4. Proof of Theorem 3.1.** Choose  $f \in L_*^2(\pi)$  such that  $\|f\|_* = 1$ . We define

$$568 \quad \mathcal{E}_i(f, f) \stackrel{\text{def}}{=} \frac{1}{2} \int_{\mathbf{B}_i} \int_{\mathbf{B}_i} (f(y) - f(x))^2 \pi_i(dx) K_i(x, dy).$$

569 From the definition

$$571 \quad (5.7) \quad 2\mathcal{E}(f, f) = \int_{\mathcal{X}} \int_{\mathcal{X}} (f(y) - f(x))^2 \pi(dx) \left[ \sum_{i \in \mathcal{I}} \pi(i|x) K_i(x, dy) \right]$$

$$572 \quad = \sum_{i \in \mathcal{I}} \pi(i) \int_{\mathcal{X}} \int_{\mathcal{X}} (f(y) - f(x))^2 \pi_i(dx) K_i(x, dy)$$

$$573 \quad \geq 2 \sum_{i \in \mathcal{I}} \pi(i) \mathcal{E}_i(f, f) \geq 2 \sum_{i \in \mathcal{I}_0} \pi(i) \mathcal{E}_i(f, f).$$

575 Using  $\bar{\mathbf{B}} = \cup_{i \in \mathcal{I}_0} \{i\} \times \mathbf{B}_i$ , and  $\bar{\mathbf{B}}^c \stackrel{\text{def}}{=} (\mathcal{I} \times \mathcal{X}) \setminus \bar{\mathbf{B}}$ , we have,

$$577 \quad (5.8) \quad 2\text{Var}_\pi(f) = \int_{\bar{\mathbf{B}}} \int_{\bar{\mathbf{B}}} (f(y) - f(x))^2 \bar{\pi}(di, dx) \bar{\pi}(dj, dy)$$

$$578 \quad + 2 \int_{\bar{\mathbf{B}}} \int_{\bar{\mathbf{B}}^c} (f(y) - f(x))^2 \bar{\pi}(di, dx) \bar{\pi}(dj, dy)$$

$$579 \quad + \int_{\bar{\mathbf{B}}^c} \int_{\bar{\mathbf{B}}^c} (f(y) - f(x))^2 \bar{\pi}(di, dx) \bar{\pi}(dj, dy).$$

581 For  $\bar{\mathbf{B}}$  as in (3.8), and expanding the first term on the right hand side of (5.8) it follows that

$$583 \quad (5.9) \quad 2 \left( \text{Var}_\pi(f) - \frac{\zeta}{2} \right) \leq \sum_{i \in \mathcal{I}_0} \pi(i)^2 \int_{\mathbf{B}_i} \int_{\mathbf{B}_i} (f(y) - f(x))^2 \pi_i(dx) \pi_i(dy)$$

$$584 \quad + \sum_{i \neq j, i, j \in \mathcal{I}_0} \pi(i) \pi(j) \pi_i(\mathbf{B}_i) \pi_j(\mathbf{B}_j) \int_{\mathbf{B}_i} \int_{\mathbf{B}_j} (f(y) - f(x))^2 \frac{\pi_i(dx)}{\pi_i(\mathbf{B}_i)} \frac{\pi_j(dy)}{\pi_j(\mathbf{B}_j)}.$$

586 Given an edge  $e$  in  $\mathcal{G}$ , let us write  $e_-$  and  $e_+$  to denote the two incident nodes of the edge.  
 587 For  $i \neq j \in \mathcal{I}_0$ , let  $\gamma_{ij}$  denotes the chosen canonical path between  $i$  and  $j$ , and let  $i_0, i_1, \dots, i_\ell$   
 588 be the nodes on that canonical path (with  $i_0 = i$ , and  $i_\ell = j$ ). By introducing generic  
 589 variables  $z_{i_k} \in \mathbf{B}_{i_k}$ , one can write  $f(z_{i_\ell}) - f(z_{i_0}) = \sum_{k=1}^{\ell} f(z_{i_k}) - f(z_{i_{k-1}})$ . Using this and the  
 590 Cauchy-Schwarz inequality, we have

$$592 \quad (5.10) \quad \int_{\mathbf{B}_i} \int_{\mathbf{B}_j} (f(y) - f(x))^2 \frac{\pi_i(dx)}{\pi_i(\mathbf{B}_i)} \frac{\pi_j(dy)}{\pi_j(\mathbf{B}_j)}$$

$$593 \quad \leq |\gamma_{ij}| \sum_{e \in \gamma_{ij}} \int_{\mathbf{B}_{e_-}} \int_{\mathbf{B}_{e_+}} (f(y) - f(x))^2 \frac{\pi_{e_-}(dx)}{\pi_{e_-}(\mathbf{B}_{e_-})} \frac{\pi_{e_+}(dy)}{\pi_{e_+}(\mathbf{B}_{e_+})},$$

594

595 where  $|\gamma_{ij}|$  denotes the number of edges on the canonical path  $\gamma_{ij}$ . By Lemma 5.1 and using  
 596 also the assumption that  $\pi_i(\mathbf{B}_i) \geq 1/2$ , the summation on the right-hand side of (5.10) is  
 597 upper bounded by

$$\begin{aligned}
 598 \quad & \frac{4}{\kappa} \sum_{e \in \gamma_{ij}} \int_{\mathbf{B}_{e_-}} \int_{\mathbf{B}_{e_-}} (f(y) - f(x))^2 \pi_{e_-}(dx) \pi_{e_-}(dy) \\
 599 \quad & + \frac{4}{\kappa} \sum_{e \in \gamma_{ij}} \int_{\mathbf{B}_{e_+}} \int_{\mathbf{B}_{e_+}} (f(y) - f(x))^2 \pi_{e_+}(dx) \pi_{e_+}(dy) \\
 600 \quad & \leq \frac{8}{\kappa} \sum_{\iota \in \gamma_{ij}} \int_{\mathbf{B}_\iota} \int_{\mathbf{B}_\iota} (f(y) - f(x))^2 \pi_\iota(dx) \pi_\iota(dy),
 \end{aligned}$$

603 where the summation  $e \in \gamma_{ij}$  is taken over all edges along the path  $\gamma_{ij}$  whereas the summation  
 604  $\iota \in \gamma_{ij}$  is taken over all nodes  $\iota$  along the path  $\gamma_{ij}$  including  $i$  and  $j$ . Hence

$$\begin{aligned}
 605 \quad & (5.11) \quad \sum_{i \neq j, i, j \in \mathfrak{l}_0} \pi(i) \pi(j) \pi_i(\mathbf{B}_i) \pi_j(\mathbf{B}_j) \int_{\mathbf{B}_i} \int_{\mathbf{B}_j} (f(y) - f(x))^2 \frac{\pi_i(dx)}{\pi_i(\mathbf{B}_i)} \frac{\pi_j(dy)}{\pi_j(\mathbf{B}_j)} \\
 606 \quad & \leq \frac{8}{\kappa} \sum_{\iota \in \mathfrak{l}_0} \pi(\iota) \int_{\mathbf{B}_\iota} \int_{\mathbf{B}_\iota} (f(y) - f(x))^2 \pi_\iota(dx) \pi_\iota(dy) \sum_{i, j \in \mathfrak{l}_0: \gamma_{ij} \ni \iota} |\gamma_{ij}| \frac{\pi(i) \pi(j)}{\pi(\iota)},
 \end{aligned}$$

609 which together with (5.9) yields

$$610 \quad (5.12) \quad 2 \left( \text{Var}_\pi(f) - \frac{\zeta}{2} \right) \leq \left( 1 + \frac{8\mathfrak{m}_1}{\kappa} \right) \sum_{i \in \mathfrak{l}_0} \pi(i) \int_{\mathbf{B}_i} \int_{\mathbf{B}_i} (f(y) - f(x))^2 \pi_i(dx) \pi_i(dy).$$

613 From the definition of  $\text{SpecGap}_i(K_i)$ , we have

$$614 \quad (5.13) \quad \int_{\mathbf{B}_i} \int_{\mathbf{B}_i} (f(y) - f(x))^2 \pi_i(dx) \pi_i(dy) \leq \frac{2\mathcal{E}_i(f, f)}{\text{SpecGap}_i(K_i)},$$

615 which we use in (5.12), to arrive at

$$616 \quad (5.14) \quad \left( \text{Var}_\pi(f) - \frac{\zeta}{2} \right) \leq \frac{\left( 1 + \frac{8\mathfrak{m}_1}{\kappa} \right)}{\min_{i \in \mathfrak{l}_0} \text{SpecGap}_i(K_i)} \sum_{i \in \mathfrak{l}_0} \pi(i) \mathcal{E}_i(f, f).$$

617 (5.14) and (5.7) together yield,

$$618 \quad \frac{\mathcal{E}(f, f)}{\left( \text{Var}_\pi(f) - \frac{\zeta}{2} \right)} \geq \frac{\min_{i \in \mathfrak{l}_0} \text{SpecGap}_i(K_i)}{1 + \frac{8\mathfrak{m}_1}{\kappa}} \geq \frac{\kappa}{1 + 8\mathfrak{m}_1} \min_{i \in \mathfrak{l}_0} \text{SpecGap}_i(K_i),$$

620 which together with the definition (2.3) implies the stated bound.  $\square$

621 **5.5. Proof of Theorem 4.2.** We start with some basic calculations on the model.

622 **Lemma 5.2.** For  $\delta, \vartheta \in \Delta$  such that  $\vartheta \supseteq \delta$ , setting  $\tau \stackrel{\text{def}}{=} \frac{1}{\sigma^2} \left( \frac{1}{\rho_1} - \frac{1}{\rho_0} \right)$ , we have

623

$$624 \quad (5.15) \quad \frac{\Pi(\vartheta|z)}{\Pi(\delta|z)} = \left( \frac{1}{p^u} \right)^{\|\vartheta\|_0 - \|\delta\|_0} \frac{e^{\frac{\tau}{2\sigma^2} z' L_\delta^{-1} X_{(\vartheta-\delta)}} \left( I_{\|\vartheta-\delta\|_0} + \tau X'_{(\vartheta-\delta)} L_\delta^{-1} X_{(\vartheta-\delta)} \right)^{-1} X'_{(\vartheta-\delta)} L_\delta^{-1} z}{\sqrt{\det \left( I_{\|\vartheta-\delta\|_0} + \tau X'_{(\vartheta-\delta)} L_\delta^{-1} X_{(\vartheta-\delta)} \right)}}.$$

625

626 *Proof.* We start with some basic calculations on the model. For any  $\vartheta, \delta \in \Delta$ , we have

$$627 \quad \frac{\Pi(\vartheta|z)}{\Pi(\delta|z)} = \frac{\omega_\vartheta}{\omega_\delta} \left( \frac{\rho_1}{\rho_0} \right)^{\frac{\|\vartheta\|_0 - \|\delta\|_0}{2}} \frac{\int_{\mathbb{R}^p} e^{-\frac{1}{2\sigma^2} \|z - Xu\|_2^2 - \frac{1}{2} u' D_{(\vartheta)}^{-1} u} \mathrm{d}u}{\int_{\mathbb{R}^p} e^{-\frac{1}{2\sigma^2} \|z - Xu\|_2^2 - \frac{1}{2} u' D_{(\delta)}^{-1} u} \mathrm{d}u}.$$

$$628 \quad = \frac{\omega_\vartheta}{\omega_\delta} \left( \frac{\rho_1}{\rho_0} \right)^{\frac{\|\vartheta\|_0 - \|\delta\|_0}{2}} \frac{\sqrt{\det \left( \sigma^2 D_{(\delta)}^{-1} + X'X \right)} e^{\frac{1}{2\sigma^2} z' X \left( \sigma^2 D_{(\vartheta)}^{-1} + X'X \right)^{-1} X'z}}{\sqrt{\det \left( \sigma^2 D_{(\vartheta)}^{-1} + X'X \right)} e^{\frac{1}{2\sigma^2} z' X \left( \sigma^2 D_{(\delta)}^{-1} + X'X \right)^{-1} X'z}}.$$

629 By the determinant lemma ( $\det(A + UV') = \det(A) \det(I_m + V'A^{-1}U)$  valid for any invertible  
630 matrix  $A \in \mathbb{R}^{n \times n}$ , and  $U, V \in \mathbb{R}^{n \times m}$ ) we have

$$631 \quad \left( \frac{\rho_1}{\rho_0} \right)^{\frac{\|\vartheta\|_0 - \|\delta\|_0}{2}} \frac{\sqrt{\det \left( \sigma^2 D_{(\delta)}^{-1} + X'X \right)}}{\sqrt{\det \left( \sigma^2 D_{(\vartheta)}^{-1} + X'X \right)}} = \sqrt{\frac{\det \left( I_n + \frac{1}{\sigma^2} X D_{(\delta)} X' \right)}{\det \left( I_n + \frac{1}{\sigma^2} X D_{(\vartheta)} X' \right)}}.$$

632 By the Woodbury identity which states that for any set of matrices  $U, V, A, C$  with matching  
633 dimensions,  $(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$ , we have

634

$$635 \quad X \left( \sigma^2 D_{(\delta)}^{-1} + X'X \right)^{-1} X' = \frac{1}{\sigma^2} X D_{(\delta)} X' - \frac{1}{\sigma^4} X D_{(\delta)} X' \left( I_n + \frac{1}{\sigma^2} X D_{(\delta)} X' \right)^{-1} X D_{(\delta)} X'$$

$$636 \quad = I_n - \left( I_n + \frac{1}{\sigma^2} X D_{(\delta)} X' \right)^{-1}.$$

637

638 so that,

$$639 \quad \frac{e^{\frac{1}{2\sigma^2} z' X \left( \sigma^2 D_{(\vartheta)}^{-1} + X'X \right)^{-1} X'z}}{e^{\frac{1}{2\sigma^2} z' X \left( \sigma^2 D_{(\delta)}^{-1} + X'X \right)^{-1} X'z}} = \frac{e^{\frac{1}{2\sigma^2} z' \left( I_n + \frac{1}{\sigma^2} X D_{(\delta)} X' \right)^{-1} z}}{e^{\frac{1}{2\sigma^2} z' \left( I_n + \frac{1}{\sigma^2} X D_{(\vartheta)} X' \right)^{-1} z}}.$$

640 We combine these developments together to conclude that

$$641 \quad (5.16) \quad \frac{\Pi(\vartheta|z)}{\Pi(\delta|z)} = \frac{\omega_\vartheta}{\omega_\delta} \sqrt{\frac{\det(L_\delta)}{\det(L_\vartheta)}} \frac{e^{\frac{1}{2\sigma^2} z' L_\delta^{-1} z}}{e^{\frac{1}{2\sigma^2} z' L_\vartheta^{-1} z}},$$

642 where, for  $\delta \in \Delta$ , we recall the definition  $L_\delta \stackrel{\text{def}}{=} I_n + \frac{1}{\sigma^2} X D_{(\delta)} X'$ . If  $\vartheta \supseteq \delta$ , setting  $\tau \stackrel{\text{def}}{=} \frac{1}{\sigma^2} \left( \frac{1}{\rho_1} - \frac{1}{\rho_0} \right) < 1/(\sigma^2 \rho_1)$ , it is easily seen that

$$644 \quad L_\vartheta = L_\delta + \tau \sum_{j: \delta_j=0, \vartheta_j=1} X_j X_j'.$$

645 The determinant lemma then gives

$$646 \quad \frac{\det(L_\vartheta)}{\det(L_\delta)} = \det \left( I_{\|\vartheta-\delta\|_0} + \tau X'_{(\vartheta-\delta)} L_\delta^{-1} X_{(\vartheta-\delta)} \right).$$

647 And the Woodbury identity gives

$$648 \quad L_\vartheta^{-1} = L_\delta^{-1} - \tau L_\delta^{-1} X_{(\vartheta-\delta)} \left( I_{\|\vartheta-\delta\|_0} + \tau X'_{(\vartheta-\delta)} L_\delta^{-1} X_{(\vartheta-\delta)} \right)^{-1} X'_{(\vartheta-\delta)} L_\delta^{-1}.$$

649 Combining the last two display in (5.16) yields the stated results. ■

650 **Lemma 5.3.** *Assume H2. Let  $\varrho$  and  $\underline{\theta}_*$  be as in Theorem 4.2. For  $z \in \mathcal{E}_0$ , we have*

651

$$652 \quad (5.17) \quad \max_{\delta \in \Delta_s} \max_{j: \delta_{*j}=0} |X_j' L_\delta^{-1} z| \leq \sqrt{\varrho n \log(p)},$$

653

654

$$\text{and} \quad \max_{\delta \in \Delta_s} \max_{j: \delta_{*j}=1} |X_j' L_\delta^{-1} z| \leq \|\theta_*\|_\infty n + \sqrt{\varrho n \log(p)}.$$

655 Furthermore, if  $n \geq 4\varrho \log(p)/(\theta_*^2 \lambda_1^2)$ , then

$$656 \quad \min_{\delta \in \Delta_s} \min_{j: \delta_{*j}=1} |X_j' L_\delta^{-1} z| \geq \frac{\lambda_1}{2} \theta_* n.$$

657 *Proof.* Set  $V \stackrel{\text{def}}{=} (z - X \theta_*)/\sigma$ , so that

$$658 \quad z = \sigma V + \sum_{k: \delta_{*k}=1} \theta_{*k} X_k,$$

659 and

$$660 \quad X_j' L_\delta^{-1} z = \sigma X_j' L_\delta^{-1} V + \sum_{k: \delta_{*k}=1} \theta_{*,k} X_j' L_\delta^{-1} X_k.$$

661 For  $z \in \mathcal{E}_0$ ,  $|X_j' L_\delta^{-1} V| \leq \sqrt{c_0 n \log(p)}$ . If  $\delta_{*j} = 0$  and  $\delta_{*k} = 1$ , then  $|X_j' L_\delta^{-1} X_k| \leq \mathcal{C}(s) \sqrt{n \log(p)}$ . ■

662 Hence

$$663 \quad \max_{\delta \in \Delta_s} \max_{j: \delta_{*j}=0} |X_j' L_\delta^{-1} z| \leq \left( \sigma \sqrt{c_0} + \mathcal{C}(s) \sum_{k: \delta_{*k}=1} |\theta_{*k}| \right) \sqrt{n \log(p)} \leq \sqrt{\varrho n \log(p)}.$$

664 If  $\delta_{*j} = 1$ , then

$$665 \quad X_j' L_\delta^{-1} z = \sigma X_j' L_\delta^{-1} V + \theta_{*,j} X_j' L_\delta^{-1} X_j + \sum_{k \neq j: \delta_{*k}=1} \theta_{*,k} X_j' L_\delta^{-1} X_k.$$

666 Since  $X_j' L_\delta^{-1} X_j \leq \|X_j\|_2^2 = n$ , this implies, as we have done above that  $|X_j' L_\delta^{-1} z| \leq \|\theta_\star\|_\infty n +$   
 667  $\sqrt{\varrho n \log(p)}$ . Similarly, if  $\delta_{\star j} = 1$ , then  $|X_j' L_\delta^{-1} z| \geq \frac{1}{2} |\theta_{\star, j}| X_j' L_\delta^{-1} X_j$ , provided that we have  
 668  $\sqrt{\varrho n \log(p)} \leq \frac{1}{2} |\theta_{\star, j}| X_j' L_\delta^{-1} X_j$ . Then using the definition of  $\lambda_1$ , we get

$$669 \quad \min_{\delta \in \Delta_s} \min_{j: \delta_{\star j}=1} |X_j' L_\delta^{-1} z| \geq \frac{\lambda_1}{2} \theta_\star n. \quad \blacksquare$$

670 *Proof of Theorem 4.2.* Fix  $\zeta_0 \in (0, 1)$ . We will apply Lemma 2.1 with  $\|\cdot\|_\star = \|\cdot\|_{\pi, \infty}$ .  
 671 Since  $Kf$  is bounded when  $f$  is bounded, the kernel  $K$  satisfies (2.2) with this choice of  
 672  $\|\cdot\|_\star$ . We recall that the initial distribution is taken as  $\nu_0 = \Pi(\cdot | \delta^{(i)}, z)$ , for some initial choice  
 673  $\delta^{(i)} \in \Delta_s$ . Let  $f_0$  be the density of  $\nu_0$  with respect to  $\Pi(\cdot | z)$ . We Lemma 2.1 with  $\zeta = 0$  to  
 674 conclude that

$$675 \quad (5.18) \quad \|\nu_0 K^N - \Pi(\cdot | z)\|_{\text{tv}}^2 \leq \zeta_0^2, \quad \text{for } N \geq \frac{1}{\text{SpecGap}_0(K)} \log \left( \frac{\text{Var}_\pi(f_0)}{\zeta_0^2} \right).$$

676 To bound the spectral gap we apply Theorem 3.1 with the choices  $\zeta = 0$ ,  $\mathsf{l} = \Delta$ ,  $\mathsf{l}_0 = \Delta_s$ ,  
 677 and  $\mathsf{B}_\delta = \mathbb{R}^p$ , and with a graph on  $\Delta_s$  constructed as follows: we put an edge between  $\delta^{(1)}$   
 678 and  $\delta^{(2)}$  if  $\delta^{(1)} \supseteq \delta^{(2)}$ , or  $\delta^{(2)} \supseteq \delta^{(1)}$ , and  $\|\delta^{(2)} - \delta^{(1)}\|_0 = 1$  (in other words the models  $\delta^{(1)}$  and  
 679  $\delta^{(2)}$  differ only in one variable). Clearly (3.8) holds, since  $\Pi(\Delta_s | z) = 1$ . We then conclude  
 680 from Theorem 3.1 that

$$681 \quad (5.19) \quad \text{SpecGap}_0(K) \geq \frac{\kappa}{1 + 8\mathfrak{m}_1}.$$

682 To bound the constants  $\kappa$  and  $\mathfrak{m}_1$  we develop a similar argument as in [Yang et al.(2016)].  
 683 Given  $\delta \in \Delta_s$ , we call  $\min(\delta, \delta_\star)$  the skeleton of  $\delta$ , and we let  $\mathcal{S} \stackrel{\text{def}}{=} \{\min(\delta, \delta_\star), \delta \in \Delta_s\}$  be the  
 684 set of all possible skeletons. Basically  $\mathcal{S}$  is the set of submodels of the true model  $\delta_\star$ . Given  
 685  $\delta \in \Delta_s$ , we build our canonical path from  $\delta$  to  $\delta_\star$  as follows. First we build a path from  $\delta$  to  
 686 its skeleton (that is  $\min(\delta, \delta_\star)$ ) by successively removing from the model  $\delta$  the variables  $X_j$  for  
 687 which  $\delta_j = 1$  and  $\delta_{\star j} = 0$ , in reverse index ordering. Then we build a path from the skeleton  
 688 to  $\delta_\star$  by adding to the skeleton the variables  $X_j$  for which  $\delta_j = 0$  and  $\delta_{\star j} = 1$  in their index  
 689 ordering. For example, if  $p = 6$ ,  $\delta_\star = (1, 1, 1, 0, 0)$  and  $\delta = (0, 0, 1, 0, 1, 1)$ , then our canonical  
 690 path from  $\delta$  to  $\delta_\star$  is

$$691 \quad (0, 0, 1, 0, 1, 1) \rightarrow (0, 0, 1, 0, 1, 0) \rightarrow (0, 0, 1, 0, 0, 0) \rightarrow (1, 0, 1, 0, 0, 0) \rightarrow (1, 1, 1, 0, 0, 0).$$

692 Given  $\delta^{(1)}, \delta^{(2)} \in \Delta_s$ , let  $\delta^{(1,2)}$  be the node where the canonical path from  $\delta^{(1)}$  to  $\delta_\star$  and the  
 693 canonical path from  $\delta^{(2)}$  to  $\delta_\star$  meet for the first time. Our canonical path  $\gamma_{\delta^{(1)}, \delta^{(2)}}$  between  $\delta^{(1)}$   
 694 and  $\delta^{(2)}$  is then defined as follows. Follow the canonical path from  $\delta^{(1)}$  towards  $\delta_\star$  until  $\delta^{(1,2)}$ ,  
 695 then reverse direction and follow the path from  $\delta^{(1,2)}$  until  $\delta^{(2)}$ . For instance if  $p = 6$ ,  $\delta_\star =$   
 696  $(1, 1, 1, 0, 0, 0)$  and  $\delta^{(1)} = (0, 1, 0, 0, 1, 1)$ , and  $\delta^{(2)} = (1, 1, 0, 1, 1, 0)$ , then  $\delta^{(1,2)} = (1, 1, 0, 0, 0, 0)$ ,  
 697 and our chosen canonical path from  $\delta^{(1)}$  to  $\delta^{(2)}$  is

$$698 \quad (0, 1, 0, 0, 1, 1) \rightarrow (0, 1, 0, 0, 1, 0) \rightarrow (0, 1, 0, 0, 0, 0) \rightarrow (1, 1, 0, 0, 0, 0) \rightarrow (1, 1, 0, 1, 0, 0) \rightarrow (1, 1, 0, 1, 1, 0).$$

699 We claim that for the canonical paths constructed above we have

$$700 \quad (5.20) \quad \mathbf{m}_1 \stackrel{\text{def}}{=} \max_{\delta \in \Delta_s} \sum_{\delta^{(1)}, \delta^{(2)} \in \Delta_s: \gamma_{\delta^{(1)}, \delta^{(2)}} \ni \delta} |\gamma_{\delta^{(1)}, \delta^{(2)}}| \frac{\pi(\delta^{(1)}|z)\pi(\delta^{(2)}|z)}{\pi(\delta|z)} \leq 8s,$$

701 and

$$702 \quad (5.21) \quad \kappa \stackrel{\text{def}}{=} \min_{\delta^{(1)} \sim \delta^{(2)}} \int_{\mathbb{R}^p} \min \left( \Pi(\theta|\delta^{(1)}, z), \Pi(\theta|\delta^{(2)}, z) \right) d\theta$$

$$704 \quad \geq \frac{1}{2} \min \left( 1, \sqrt{\frac{\sigma^2 \rho_0}{2n}} \right) \min \left( 1, e^{\frac{1}{4\sigma^2}(n\lambda_1 - 8\sigma^2 \rho_0)} \right) p^{-\frac{2\rho_0}{n} \frac{e}{\lambda^2}}.$$

706 where the minimum is taken over all connected pairs of nodes  $\delta^{(1)}, \delta^{(2)}$ . Furthermore, we claim  
707 that we can bound the variance of the initial density and get

$$708 \quad (5.22) \quad \log \left( \frac{\text{Var}_\pi(f_0)}{\zeta_0^2} \right) \leq A \left( \log \left( \frac{1}{\zeta_0} \right) + \frac{su(1 + \|\theta_\star\|_\infty^2)n}{\sigma^2 \lambda} \right),$$

709 for some absolute constant  $A$ . (5.20) and (5.21) shows that

$$710 \quad (5.23) \quad \text{SpecGap}_0(K) \geq \frac{A}{s} \min \left( 1, \sqrt{\frac{\sigma^2 \rho_0}{2n}} \right) \min \left( 1, e^{\frac{1}{4\sigma^2}(n\lambda_1 - 8\sigma^2 \rho_0)} \right) p^{-\frac{2\rho_0}{n} \frac{e}{\lambda^2}}$$

711 for some absolute constant  $A$ . We put (5.23) together with (5.22) and (5.18) to reach the  
712 stated conclusion. The remaining of the proof consists in establishing the claims (5.20), (5.21)  
713 and (5.22).

714 **Proof of Equation (5.20).** For  $\delta^{(1)}, \delta^{(2)} \in \Delta_s$ , we will use the obvious bound

$$715 \quad |\gamma_{\delta^{(1)}, \delta^{(2)}}| \leq 2s.$$

716 Given  $\delta \in \Delta_s$ , we denote  $\Lambda(\delta)$  the set of all  $\delta^{(1)} \in \Delta_s$  such that the canonical path from  $\delta^{(1)}$   
717 to  $\delta_\star$  goes through  $\delta$ . Using this we can bound  $\mathbf{m}_1$  as

$$718 \quad (5.24) \quad \mathbf{m}_1 \leq 2s \max_{\delta \in \Delta_s} \sum_{\delta^{(1)} \in \Lambda(\delta)} \sum_{\delta^{(2)} \in \Delta_s} \frac{\pi(\delta^{(1)}|z)\pi(\delta^{(2)}|z)}{\pi(\delta|z)} \leq 2s \max_{\delta \in \Delta_s} \sum_{\delta^{(1)} \in \Lambda(\delta)} \frac{\pi(\delta^{(1)}|z)}{\pi(\delta|z)}.$$

721 Let  $\mathcal{S} \stackrel{\text{def}}{=} \{\min(\delta, \delta_\star), \delta \in \Delta_s\}$  be the set of all possible skeletons. Take  $\delta^{(1)} \in \Lambda(\delta)$ . We  
722 will distinguish whether  $\delta \in \mathcal{S}$  or not. Suppose  $\delta \notin \mathcal{S}$ . Therefore, traveling the canonical  
723 path from  $\delta^{(1)}$  toward  $\delta_\star$  we arrive at  $\delta$  by removing only non-significant variables. Therefore,  
724 assuming that  $\|\delta^{(1)}\|_0 = \|\delta\|_0 + \ell$ , and using (5.15), and H2, we have

$$725 \quad (5.25) \quad \frac{\pi(\delta^{(1)}|z)}{\pi(\delta|z)} \leq \frac{1}{p^{u\ell}} \exp \left( \frac{\tau}{2\sigma^2(1+n\tau\lambda)} \sum_{j: \delta_j^{(1)}=1, \delta_j=0} (X'_j L_\delta^{-1} z)^2 \right) \leq \frac{e^{\frac{\ell \bar{Q}_0}{n\lambda}}}{p^{u\ell}},$$

726 where  $\bar{Q}_0 = \max_{j: \delta_j^{(1)}=1, \delta_j=0} (X'_j L_\delta^{-1} z)^2$ . From Lemma 5.3, we get  $\bar{Q}_0 \leq \varrho n \log(p)$ . Using this  
 727 and the trivial inequality  $\binom{p}{\ell} \leq p^\ell$ , it follows that

$$728 \quad \sum_{\delta^{(1)} \in \Lambda(\delta)} \frac{\pi(\delta^{(1)}|z)}{\pi(\delta|z)} \leq \sum_{\ell=0}^{s-\|\delta\|_0} \sum_{\delta^{(1)} \in \Lambda(\delta): \|\delta^{(1)}\|_0 = \|\delta\|_0 + \ell} \frac{\pi(\delta^{(1)}|z)}{\pi(\delta|z)} \leq \sum_{\ell=0}^s \left( \frac{p^{\frac{\varrho}{2\sigma^2\lambda}}}{p^{u-1}} \right)^\ell \leq 2,$$

731 under the assumption that  $\sigma^2 u \lambda \geq \varrho$ , and  $u > 4$ . Suppose now that  $\delta \in \mathcal{S}$ . Then  $\Lambda(\delta)$  is  
 732 comprised of the elements of  $\Delta_s$  whose skeletons are subsets of  $\delta$ . Hence

$$733 \quad \sum_{\delta^{(1)} \in \Lambda(\delta)} \frac{\pi(\delta^{(1)}|z)}{\pi(\delta|z)} = \sum_{\delta_0 \in \mathcal{S}: \delta \supseteq \delta_0} \frac{\pi(\delta_0|z)}{\pi(\delta|z)} \sum_{\delta^{(1)} \in \Lambda(\delta): \min(\delta^{(1)}, \delta_*) = \delta_0} \frac{\pi(\delta^{(1)}|z)}{\pi(\delta_0|z)}.$$

734 The inner summation can be upper bounded by 2 as above. If  $\delta \supseteq \delta_0$  and  $\|\delta\|_0 = \|\delta_0\|_0 + r$ ,  
 735 we apply (5.15) again and get,

$$736 \quad \frac{\pi(\delta_0|z)}{\pi(\delta|z)} \leq \left( p^u \sqrt{1 + \frac{ns_*}{\sigma^2 \rho_1}} e^{-\frac{\tau \bar{Q}_3}{2\sigma^2(1+\tau s_* n)}} \right)^r \leq \left( p^{u+a} e^{-\frac{\bar{Q}_3}{4\sigma^2 s_* n}} \right)^r,$$

737 where we use H3-(3) to obtain  $\tau/(1 + \tau s_* n) \geq 1/(2s_* n)$ , and  $\sqrt{1 + \frac{ns_*}{\sigma^2 \rho_1}} \leq p^a$ , and where  
 738  $\bar{Q}_3 \stackrel{\text{def}}{=} \min_{j: \delta_{0j}=0, \delta_{*j}=1} (X'_j L_{\delta_0}^{-1} z)^2$ . From Lemma 5.3 we get  $\bar{Q}_3 \geq \frac{\varrho^2}{4} \lambda_1^2 n^2$ , under the sample  
 739 condition  $n \geq 4\varrho \log(p)/(\varrho_*^2 \lambda_1^2)$  which is implied by (4.10). We conclude that

$$740 \quad \max_{\delta \in \mathcal{S}} \sum_{\delta^{(1)} \in \Lambda(\delta)} \frac{\pi(\delta^{(1)}|z)}{\pi(\delta|z)} \leq 2 \sum_{\delta_0 \in \mathcal{S}: \delta \supseteq \delta_0} \frac{\pi(\delta_0|z)}{\pi(\delta|z)} \leq 2 \sum_{r=0}^{s_*} s_*^r \left( p^{u+a} e^{-\frac{\varrho_*^2 \lambda_1^2 n}{4\sigma^2 s_*}} \right)^r \leq 4,$$

741 using the sample size condition in (4.10). This proves the claim (5.20).

742 **Proof of Equation (5.21).** Fix  $\delta^{(1)}, \delta^{(2)} \in \Delta_s$ , such that  $\delta^{(1)} \supseteq \delta^{(2)}$ , or  $\delta^{(2)} \subseteq \delta^{(1)}$ , and  
 743  $\|\delta^{(2)} - \delta^{(1)}\|_0 = 1$ . Without any loss of generality, suppose that  $\delta^{(2)} \supseteq \delta^{(1)}$ , and their difference  
 744 occurs on component  $j$ :  $\delta_j^{(2)} = 1$ , while  $\delta_j^{(1)} = 0$ . Then for all  $\theta \in \mathbb{R}^p$ , we have

$$745 \quad \frac{\Pi(\theta|\delta^{(1)}, z)}{\Pi(\theta|\delta^{(2)}, z)} = \left( \frac{\int_{\mathbb{R}^p} e^{-\frac{1}{2\sigma^2} \|z - X\theta\|_2^2 - \frac{1}{2} \theta' D_{(\delta^{(2)})}^{-1} \theta} d\theta}{\int_{\mathbb{R}^p} e^{-\frac{1}{2\sigma^2} \|z - X\theta\|_2^2 - \frac{1}{2} \theta' D_{(\delta^{(1)})}^{-1} \theta} d\theta} \right) e^{-(\rho_0 - \rho_1) \frac{\theta_j^2}{2}}.$$

746 Let  $A$  denote the ratio of integrals in the last display. We can then write

$$747 \quad \int_{\mathbb{R}^p} \min \left( \Pi(\theta|\delta^{(1)}, z), \Pi(\theta|\delta^{(2)}, z) \right) d\theta = \int_{\mathbb{R}} \min \left( 1, A e^{-(\rho_0 - \rho_1) \frac{\theta_j^2}{2}} \right) \Pi(\theta_j|\delta^{(2)}, z) d\theta_j.$$

748 Recall from (4.3) that the  $j$  marginal under  $\Pi(\theta_j|\delta^{(2)}, z)$  is the Gaussian distribution  $\mathbf{N}(\mu_j, \sigma_j^2)$ ,  
 749 where

$$750 \quad \sigma_j = \sigma \sqrt{e'_j \Sigma_{\delta^{(2)}} e_j}, \quad \text{and} \quad \mu_j = e'_j \Sigma_{\delta^{(2)}} X' z, \quad 1 \leq j \leq p,$$

751 and where  $e_j$  denotes the  $j$ -th unit vector. Hence, for  $Z \sim \mathbf{N}(0, 1)$ ,

$$752 \quad (5.26) \quad \int_{\mathbb{R}^p} \min \left( \Pi(\theta|\delta^{(1)}, z), \Pi(\theta|\delta^{(2)}, z) \right) d\theta = \mathbb{E} \left[ \min \left( 1, A e^{-\frac{(\rho_0 - \rho_1)}{2}(\mu_j + \sigma_j Z)^2} \right) \right]$$

$$754 \quad \geq \frac{1}{2} \min \left( 1, A e^{-\frac{(\rho_0 - \rho_1)}{2}(|\mu_j| + \sigma_j)^2} \right) \geq \frac{1}{2} \min \left( 1, A e^{-\rho_0(\mu_j^2 + \sigma_j^2)} \right),$$

756 using the fact that for any nonnegative function  $f$ ,  $\mathbb{E}(f(Z)) \geq \mathbb{P}(|Z| \leq 1) \min_{z: |z| \leq 1} f(z)$ . By  
757 matrix block inversion, we work out  $\sigma_j^2$  to

$$758 \quad (5.27) \quad \sigma_j^2 = \frac{\sigma^2}{\sigma^2 \rho_1 + X_j' \left( I_n + \frac{1}{\sigma^2} X_{-j} D_{(\delta^{(2)}, j)} X_{-j}' \right)^{-1} X_j} = \frac{\sigma^2}{\sigma^2 \rho_1 + X_j' L_{\delta_{-j}^{(1)}}^{-1} X_j} \leq \frac{\sigma^2}{\lambda_1 n},$$

759 where  $D_{(\delta^{(2)}, j)} = D_{(\delta^{(1)}, j)}$  is the  $(p-1)$ -dimensional matrix obtained by removing the  $j$ -th row  
760 and the  $j$ -th column of  $D_{(\delta^{(2)})}$ , and  $L_{\delta_{-j}^{(1)}} = I_n + \frac{1}{\sigma^2} X_{-j} D_{(\delta^{(1)}, j)} X_{-j}'$ . By block inversion the  
761 mean  $\mu_j$  can be written as

$$762 \quad (5.28) \quad \mu_j = e_1 \left( \begin{array}{cc} X_j' X_j + \sigma^2 \rho_1 & X_j' X_{-j} \\ X_{-j}' X_j & X_{-j}' X_{-j} + \sigma^2 D_{(\delta^{(2)}, j)}^{-1} \end{array} \right)^{-1} \begin{pmatrix} X_j' z \\ X_{-j}' z \end{pmatrix} = \frac{X_j' L_{\delta_{-j}^{(1)}}^{-1} z}{\sigma^2 \rho_1 + X_j' L_{\delta_{-j}^{(1)}}^{-1} X_j}.$$

764

765 Consider first the case where  $j$  is such that  $\delta_{\star, j} = 0$ . Note that  $X_j' L_{\delta_{-j}^{(1)}}^{-1} X_j \geq X_j' L_{\delta^{(1)}}^{-1} X_j \geq$   
766  $n \lambda_1$ . Therefore, and using Lemma 5.3, and  $z \in \mathcal{E}_0$ , we obtain

$$767 \quad |\mu_j| \leq \frac{1}{n \lambda_1} \sqrt{\varrho n \log(p)} = \frac{1}{\lambda_1} \sqrt{\frac{\varrho \log(p)}{n}}.$$

768 Consider now the case where  $\delta_{\star, j} = 1$ . Then we have

$$769 \quad \mu_j^2 \leq \frac{1}{(X_j' L_{\delta_{-j}^{(1)}}^{-1} X_j)^2} \left( \theta_{\star j} X_j' L_{\delta_{-j}^{(1)}}^{-1} X_j + \sigma \sqrt{c_0 n \log(p)} + \sum_{k: \delta_{\star k} = 1} \theta_{\star k} X_j' L_{\delta_{-j}^{(1)}}^{-1} X_k \right)^2$$

$$770 \quad \leq \frac{2}{(X_j' L_{\delta_{-j}^{(1)}}^{-1} X_j)^2} \left( \theta_{\star j}^2 (X_j' L_{\delta_{-j}^{(1)}}^{-1} X_j)^2 + \varrho n \log(p) \right)$$

$$771 \quad \leq 2\theta_{\star j}^2 + \frac{\varrho \log(p)}{\lambda_1^2 n}.$$

772 On the other hand, using (5.15), the ratio of integrals  $A$  gives

$$773 \quad A = \sqrt{\frac{\rho_0}{\rho_1}} \frac{1}{\sqrt{1 + \tau X_j' L_{\delta^{(1)}}^{-1} X_j}} \exp \left( \frac{1}{2\sigma^2} \frac{\tau (X_j' L_{\delta^{(1)}}^{-1} z)^2}{1 + \tau X_j' L_{\delta^{(1)}}^{-1} X_j} \right),$$

774 where we recall that  $\tau = (\rho_1^{-1} - \rho_0^{-1})/\sigma^2$ . Note that if  $\delta_{\star j} = 0$ , the term inside the exponential  
 775 in this last expression of  $A$  grows like  $\varrho \log(p)/\lambda_1$  which is not fast enough to face off with the  
 776 term  $-\rho_0(\mu_j^2 + \sigma_j^2)$ . Hence we use instead the trivial lower bound  $A \geq 1$  together with the  
 777 upper bounds on  $\mu_j$  and  $\sigma_j^2$  obtained above and (5.26) to conclude that

$$778 \quad (5.29) \quad \int_{\mathbb{R}^p} \min \left( \Pi(\theta|\delta^{(1)}, z), \Pi(\theta|\delta^{(2)}, z) \right) d\theta \geq \frac{1}{2} e^{-\rho_0(\mu_j^2 + \sigma_j^2)} \geq \frac{1}{2} \exp \left( -\frac{2\rho_0}{n} \frac{\varrho \log(p)}{\lambda_1^2} \right).$$

779 However if  $\delta_{\star j} = 1$ , By Lemma 5.3, and under the sample size condition (4.10) we have

$$780 \quad (X'_j L_{\delta^{(1)}}^{-1} z)^2 \geq \frac{\theta_{\star j}^2}{2} (X'_j L_{\delta^{(1)}}^{-1} X_j)^2.$$

781 Noting that  $1 \leq \tau X'_j L_{\delta^{(1)}}^{-1} X_j$ , we deduce that

$$782 \quad A \geq \sqrt{\frac{\rho_0}{\rho_1}} \frac{1}{\sqrt{1 + \frac{n}{\sigma^2 \rho_1}}} e^{\frac{\theta_{\star j}^2 X'_j L_{\delta^{(1)}}^{-1} X_j}{4\sigma^2}}.$$

783 It follows in this case that

$$785 \quad (5.30) \quad \int_{\mathbb{R}^p} \min \left( \Pi(\theta|\delta^{(1)}, z), \Pi(\theta|\delta^{(2)}, z) \right) d\theta \geq \frac{1}{2} \min \left( 1, A e^{-\rho_0 \left( 2\theta_{\star j}^2 + \frac{\varrho \log(p)}{\lambda^2 n} + \frac{\sigma^2}{\lambda n} \right)} \right)$$

$$786 \quad \geq \frac{1}{2} \min \left( 1, \sqrt{\frac{\sigma^2 \rho_0}{\sigma^2 \rho_1 + n}} \right) \min \left( 1, e^{\frac{\theta_{\star j}^2 (X'_j L_{\delta^{(1)}}^{-1} X_j - 8\sigma^2 \rho_0)}{4\sigma^2}} e^{-2\rho_0 \frac{\varrho \log(p)}{\lambda^2 n}} \right)$$

$$787 \quad \geq \frac{1}{2} \min \left( 1, \sqrt{\frac{\sigma^2 \rho_0}{2n}} \right) \min \left( 1, e^{\frac{1}{4\sigma^2} (n\lambda_1 - 8\sigma^2 \rho_0)} \right) p^{-\frac{2\rho_0}{n} \frac{\varrho}{\lambda^2}},$$

788

789 where we have used the fact that  $\min(1, ab) \geq \min(1, a) \min(1, b)$  valid for all nonnegative  
 790 numbers  $a, b, c$ . We combine (5.29) and (5.30) to obtain (5.21).

791 **Proof of Equation (5.22).** Since  $\Pi(\theta|z) = \sum_{\vartheta} \Pi(\vartheta|z) \Pi(\theta|\vartheta, z) \geq \Pi(\delta^{(i)}|z) \Pi(\theta|\delta^{(i)}, z)$ ,  
 792 we have

$$793 \quad f_0(\theta) = \frac{\Pi(\theta|\delta^{(i)}, z)}{\Pi(\theta|z)} \leq \frac{1}{\Pi(\delta^{(i)}|z)} = \frac{1}{\Pi(\delta_{\star}|z)} \frac{\Pi(\delta_{\star}|z)}{\Pi(\delta_{\star}^{(i)}|z)} \frac{\Pi(\delta^{(i)}|z)}{\Pi(\delta^{(i)}|z)},$$

794 where  $\delta_{\star}^{(i)} \stackrel{\text{def}}{=} \min(\delta^{(i)}, \delta_{\star})$ . We apply (5.15) twice (to each ratio), and use H2, to get

$$796 \quad \frac{\Pi(\delta_{\star}|z)}{\Pi(\delta_{\star}^{(i)}|z)} \frac{\Pi(\delta_{\star}^{(i)}|z)}{\Pi(\delta^{(i)}|z)} \leq p^{u(\|\delta\|_0 - \|\delta_{\star}\|_0)} \sqrt{\det \left( I_{\|\delta\|_0 - \|\delta_{\star}^{(i)}\|_0} + \tau X'_{(\delta - \delta_{\star}^{(i)})} L_{\delta_{\star}^{(i)}}^{-1} X_{(\delta - \delta_{\star}^{(i)})} \right)}$$

$$797 \quad \times e^{\frac{\tau}{2\sigma^2} z' L_{\delta_{\star}^{(i)}}^{-1} X_{(\delta_{\star} - \delta_{\star}^{(i)})} \left( I_{\|\delta_{\star} - \delta_{\star}^{(i)}\|_0} + \tau X'_{(\delta_{\star} - \delta_{\star}^{(i)})} L_{\delta_{\star}^{(i)}}^{-1} X_{(\delta_{\star} - \delta_{\star}^{(i)})} \right)^{-1} X'_{(\delta_{\star} - \delta_{\star}^{(i)})} L_{\delta_{\star}^{(i)}}^{-1} z}$$

$$798 \quad \leq p^{u(\|\delta\|_0 - s_{\star})} \left( 1 + \frac{n\|\delta\|_0}{\sigma^2 \rho_1} \right)^{\frac{\|\delta\|_0}{2}} e^{\frac{1}{2\sigma^2 n \lambda} \|X'_{(\delta_{\star} - \delta_{\star}^{(i)})} L_{\delta_{\star}^{(i)}}^{-1} z\|_2^2}.$$

799

800 Under the assumption  $p^{us_*}\Pi(\delta_*|z) \geq 1$  (H3-(1)), and since  $\|\delta\|_0 \leq s$ , we conclude that

$$801 \quad (5.31) \quad \|f_0\|_{\pi, \infty} \leq p^{us} \left(1 + \frac{ns}{\sigma^2 \rho_1}\right)^{\frac{s}{2}} e^{\frac{s_* \bar{Q}_1}{2\sigma^2 n \lambda}} \leq p^{(u+a)s} e^{\frac{s_* \bar{Q}_1}{2\sigma^2 n \lambda}},$$

802 where the second inequality uses (4.8), and where  $\bar{Q}_1 = \max_{j: \delta_{*,j}=1} (X'_j L_{\delta_*^{(j)}}^{-1} z)^2$ . From Lemma  
803 5.3, we get  $\bar{Q}_1 \leq 4n^2 \|\theta_*\|_{\infty}^2$ , using the sample size condition (4.10). (5.31) then becomes

$$804 \quad \sqrt{\text{Var}_{\pi}(f_0)} \leq \|f_0\|_{\pi, \infty} \leq p^{(u+a)s} e^{\frac{2s_* \|\theta_*\|_{\infty}^2 n}{\sigma^2 \lambda}} \leq e^{\frac{As(1+\|\theta_*\|_{\infty}^2)n}{\sigma^2 \lambda}},$$

805 for some absolute constant  $A$ . The claim follows by taking the log. ■

806 **5.6. Proof of Theorem 4.3.** The proof is very similar to the proof of Theorem 4.2.  
807 Fix  $\zeta_0 \in (0, 1)$ , and  $z \in \mathcal{E}$ . First we bound the uniform norm of the density of the initial  
808 distribution  $\nu_0$  as in (5.31). Noting here that the skeleton of  $\delta^{(i)}$  is  $\delta_*$ , we get the simpler  
809 bound

$$810 \quad \|f_0\|_{\pi, \infty} \leq 2 \left( p^u \sqrt{1 + \frac{n\text{FP}}{\sigma^2 \rho_1}} \right)^{\text{FP}} \leq 2p^{(u+a)\text{FP}}.$$

811 In view of this bound, we set

$$812 \quad (5.32) \quad \zeta = \frac{\zeta_0^2}{8} p^{-2(u+a)\text{FP}},$$

813 which gives  $\zeta \|f_0\|_{\pi, \infty}^2 \leq \zeta_0^2/2$ . Therefore, we can readily apply Lemma 2.1 with this particular  
814 value of  $\zeta$  to get

$$815 \quad (5.33) \quad \|\nu_0 K^N - \Pi(\cdot|z)\|_{\text{tv}}^2 \leq \zeta_0^2, \quad \text{for } N \geq \frac{1}{\text{SpecGap}_{\zeta}(K)} \log \left( \frac{1}{\zeta} \right).$$

816 We lower bound the approximate spectral gap via Theorem 3.1, and using the same approach  
817 as in Theorem 4.2. We apply Theorem 3.1 with the choices  $\mathfrak{l} = \Delta$ ,  $\mathfrak{l}_0 = \mathcal{D}_k$  endowed with the  
818 same graph as in proof of Theorem 4.2, and  $\mathfrak{B}_{\delta} = \mathbb{R}^p$ . First we need to check (3.8). For  $z \in \mathcal{E}$ ,  
819  $\zeta$  as in (5.32), we have

$$820 \quad \frac{10}{\zeta} (1 - \Pi(\mathcal{D}_k|z)) \leq \frac{80}{\zeta_0^2} p^{2(u+a)\text{FP}} \frac{1}{p^{\frac{u(k+1)}{2}}} \leq 1,$$

821 where the last inequality follows from condition (4.12). In other words we have  $\Pi(\mathcal{D}_k|z) \geq$   
822  $1 - (\zeta/10)$ , which by Remark 3.2 implies (3.8). We then conclude from Theorem 3.1 that

$$823 \quad (5.34) \quad \text{SpecGap}_{\zeta}(K) \geq \frac{\kappa}{1 + 8m_1},$$

824 where  $\kappa$  and  $m_1$  are defined using  $\mathcal{D}_k$ . We bound these terms as in Theorem 4.2 with some  
825 important simplifications due the facts that all models here belong to  $\mathcal{D}_k$ . In particular, since  
826  $\mathcal{D}_k \subseteq \Delta_s$ , we readily have

$$827 \quad (5.35) \quad m_1 \leq 8k.$$

828 Similarly, the lower bound on  $\kappa$  also simplifies. Because  $\delta^{(1)}$  and  $\delta^{(2)}$  can differ only at a  
 829 component  $j$  such that  $\delta_{*j} = 0$  (a non-important variable), we see that only the lower bound  
 830 (5.29) applies. Hence  $\kappa$  can be taken as

$$831 \quad (5.36) \quad \kappa = \frac{1}{2^p} \frac{-2\rho_0}{n} \frac{\rho}{\lambda_1^2}.$$

832 The theorem follows from the same calculations as in the proof of Theorem 4.2. □

834 **Appendix A. Some technical results.** We make use of the following standard Gaussian  
 835 deviation bound.

836 **Lemma A.1.** *Let  $Z \sim \mathbf{N}(0, I_m)$ , and  $u_1, \dots, u_N$  be vectors of  $\mathbb{R}^m$ . Then for all  $x \geq 0$ ,*

$$837 \quad \mathbb{P} \left[ \max_{1 \leq j \leq N} |\langle u_j, Z \rangle| > \max_{1 \leq j \leq N} \|u_j\|_2 \sqrt{2(x + \log(N))} \right] \leq \frac{2}{e^x}.$$

838 **Lemma A.2.** *Suppose that  $X \in \mathbb{R}^{n \times p}$  is a random matrix with i.i.d. standard Normal  
 839 entries. Given an integer  $s$ , and positive constants  $\sigma, \gamma$  and  $\rho$ , set*

$$840 \quad \mathcal{C}_0 \stackrel{\text{def}}{=} \max_{\delta \in \Delta: \|\delta\|_0 \leq s} \max_{i \neq j, \delta_j = 0} \left| X'_j \left( I_n + \frac{1}{\sigma^2 \rho_1} X_\delta X'_\delta + \frac{1}{\sigma^2 \rho_0} X_{\delta^c} X'_{\delta^c} \right) X_i \right|.$$

841 *Then there exist some universal finite constants  $c_0, a, A$  such that for  $n \geq As^2 \log(p)$ , the  
 842 following two statements hold with probability at least  $1 - \frac{a}{p}$ : for  $\rho_0^{-1} > 0$  taken small enough  
 843 and*

$$844 \quad (\text{A.1}) \quad \sigma^2 s \rho_1 \leq c_0 \sqrt{n \log(p)},$$

845 *it holds that*

846

$$847 \quad (\text{A.2}) \quad \mathcal{C}_0 \leq 2c_0 \sqrt{n \log(p)}, \quad \text{and}$$

$$848 \quad \min_{\delta: \|\delta\|_0 \leq s} \inf \left\{ \frac{u'(X'_{\delta^c} L_\delta^{-1} X_{\delta^c}) u}{n \|u\|_2^2}, u \in \mathbb{R}^{p-s}, 0 < \|\text{supp}(u)\|_0 \leq s \right\} \geq \frac{1}{32}.$$

849

850 *Proof.* For a matrix  $M \in \mathbb{R}^{n \times p}$  we set

$$851 \quad v(M, s) \stackrel{\text{def}}{=} \inf \left\{ \frac{u'(M'M)u}{n \|u\|_2^2} \mid u \neq 0, \|u\|_0 \leq s \right\},$$

852 and for  $\kappa_0 = 1/64$  and  $c_0 = 8$ , we define

853

$$854 \quad \mathcal{E} \stackrel{\text{def}}{=} \left\{ M \in \mathbb{R}^{n \times p} : v(M, s) \geq \kappa_0, \max_{1 \leq j \leq p} \|M_j\|_2 \leq 2\sqrt{n}, \right.$$

$$855 \quad \left. \min_{1 \leq j \leq p} \|M_j\|_2 \geq \sqrt{\frac{n}{2}}, \text{ and } \max_{j \neq k} |\langle M_j, M_k \rangle| \leq c_0 \sqrt{n \log(p)} \right\}.$$

856

857 By Theorem 1 of [Raskutti et al.(2010)], Lemma 1-(4.2) of [Laurent and Massart(2000)], and  
 858 standard Gaussian deviation bounds, we can find universal constants  $a, A$ , such that for  $n \geq$   
 859  $As \log(p)$ , we have  $\mathbb{P}(X \notin \mathcal{E}) \leq \frac{a}{p}$ . So to obtained the statement of the lemma, it suffices to  
 860 consider some arbitrary element  $X \in \mathcal{E}$  and show that (A.2) holds.

861 Fix  $\delta \in \Delta$  such that  $\|\delta\|_0 \leq s$ . We set  $M_\delta \stackrel{\text{def}}{=} I_n + \frac{1}{\sigma^2 \rho_1} X_\delta X'_\delta$ , so that  $L_\delta = M_\delta + \frac{1}{\sigma^2 \rho_0} X_{\delta^c} X'_{\delta^c}$ .  
 862 The Woodbury identity gives  
 (A.3)

$$863 \quad X'_j L_\delta^{-1} X_k = X'_j M_\delta^{-1} X_k - \frac{1}{\sigma^2 \rho_0} X'_j M_\delta^{-1} X_{\delta^c} \left( I_{\|\delta^c\|_0} + \frac{1}{\sigma^2 \rho_0} X'_{\delta^c} M_\delta^{-1} X_{\delta^c} \right)^{-1} X'_{\delta^c} M_\delta^{-1} X_k.$$

864 If  $C_1 = \max_\ell X'_\ell M_\delta^{-1} X_\ell$ , and  $C_0 = \max_{\ell \neq j, \delta_j=0} |X'_j M_\delta^{-1} X_\ell|$ , then we deduce easily from (A.3)  
 865 that for all  $j \neq k$  such that  $\delta_j = 0$ ,

$$866 \quad (\text{A.4}) \quad |X'_j L_\delta^{-1} X_k| \leq C_0 + \frac{1}{\sigma^2 \rho_0} (C_1^2 + p C_0^2).$$

867 In order to proceed, we need to bound the term  $X_j M_\delta^{-1} X_k$ . Easily, for  $X \in \mathcal{E}$ , we have

$$868 \quad X'_j M_\delta^{-1} X_j \leq \|X_j\|_2^2 \leq 4n.$$

869 Another application of the Woodbury identity gives

$$870 \quad (\text{A.5}) \quad M_\delta^{-1} = I_n - \frac{1}{\sigma^2 \rho_1} X_\delta \left( I_{\|\delta\|_0} + \frac{1}{\sigma^2 \rho} X'_\delta X_\delta \right)^{-1} X'_\delta.$$

871 Therefore, for  $k \neq j$

$$872 \quad X'_j M_\delta^{-1} X_k = X'_j X_k - \frac{1}{\sigma^2 \rho_1} X'_j X_\delta \left( I_{\|\delta\|_0} + \frac{1}{\sigma^2 \rho} X'_\delta X_\delta \right)^{-1} X'_\delta X_k.$$

873 Using  $X \in \mathcal{E}$ , we educe for  $j \neq k$ , and  $\delta_j = 0$ ,

874

$$875 \quad \frac{1}{\sigma^2 \rho_1} \left| X'_j X_\delta \left( I_{\|\delta\|_0} + \frac{1}{\sigma^2 \rho} X'_\delta X_\delta \right)^{-1} X'_\delta X_k \right| \leq \frac{1}{\kappa_0 n} \|X'_\delta X_k\|_2 \|X'_\delta X_j\|_2$$

$$876 \quad \leq \frac{c_0^2 s \log(p) + c_0 \sqrt{s \log(p)}}{\kappa_0} \leq c_0 \sqrt{n \log(p)},$$

877

878 for  $n \geq As^2 \log(p)$ , for some constant  $A$ . It follows that

$$879 \quad |X'_j M_\delta^{-1} X_k| \leq 2c_0 \sqrt{n \log(p)}.$$

880 We combine this with (A.4) to obtain that for  $j \neq k$  such that  $\delta_j = 0$ ,

881

$$882 \quad (\text{A.6}) \quad |X'_j L_\delta^{-1} X_k| \leq 3c_0 \sqrt{n \log(p)} \left( 1 + \frac{1}{\sigma^2 \rho_0} p c_0 \sqrt{n \log(p)} \right) + 16 \frac{1}{\sigma^2 \rho_0} n^2 \leq 8c_0 \sqrt{n \log(p)},$$

883

884 for  $\rho_0$  large enough. (A.6) says that  $\mathcal{C}_0 \leq 8c_0\sqrt{n\log(p)}$ , for  $X \in \mathcal{E}$ , as claimed.

885 For  $j$  such that  $\delta_j = 0$ , (A.5) gives

$$\begin{aligned}
 886 \quad X'_j M_\delta^{-1} X_j &= \|X_j\|_2^2 - \frac{1}{\sigma^2 \rho_1} X'_j X_\delta \left( I_{\|\delta\|_0} + \frac{1}{\sigma^2 \rho_1} X'_\delta X_\delta \right)^{-1} X'_\delta X_j \\
 887 \quad &\geq \|X_j\|_2^2 - \frac{\|X'_\delta X_j\|_2^2}{n\kappa_0} \\
 888 \quad (A.7) \quad &\geq \frac{n}{4},
 \end{aligned}$$

889 since  $n \geq As \log(p)$ , and by taking  $A$  large enough ( $A \geq 4c_0^2/\kappa_0$ ). Equation (??) then yields

890

$$\begin{aligned}
 891 \quad X'_j L_\delta^{-1} X_j &\geq X'_j M_\delta^{-1} X_j - \frac{1}{\sigma^2 \rho_0} \|X'_{\delta^c} M_\delta^{-1} X_j\|_2^2 \\
 892 \quad &= X'_j M_\delta^{-1} X_j - \frac{1}{\sigma^2 \rho_0} \left[ (X'_j M_\delta^{-1} X_j)^2 + \sum_{k: \delta_k=0, k \neq j} (X'_j M_\delta^{-1} X_k)^2 \right]. \\
 893
 \end{aligned}$$

894 For  $2\rho_0^{-1} \leq \sigma^2$ , it follows that

$$895 \quad X'_j L_\delta^{-1} X_j \geq \frac{n}{8} - \frac{1}{\sigma^2 \rho_0} (p - \|\delta\|_0) (4c_0^2 n \log(p)),$$

896 which together with (A.6) and (A.1) implies that for any  $u \in \mathbb{R}^p$  such that  $\delta^c \supseteq \text{supp}(u)$ , and  
897  $\|\text{supp}(u)\|_0 \leq s$ , we have

$$898 \quad u' X'_{\delta^c} L_\delta^{-1} X_{\delta^c} u \geq \frac{n}{32} \|u\|_2^2, \quad \blacksquare$$

899 as claimed.

900 **Acknowledgements.** I'm grateful to Joonha Park for pointing out a mistake in an initial  
901 draft of the manuscript.

902

## REFERENCES

- 903 [Atchade and Bhattacharyya(2018)] ATCHADE, Y. and BHATTACHARYYA, A. (2018). An approach to large-  
904 scale Quasi-Bayesian inference with spike-and-slab priors. *arXiv e-prints* arXiv:1803.10282.
- 905 [Bakry et al.(2008)] BAKRY, D., BARTHE, F., CATTIAUX, P. and GUILLIN, A. (2008). A simple proof of the  
906 poincaré inequality for a large class of probability measures including the log-concave case. *Electronic*  
907 *Communications in Probability* **13**.
- 908 [Bakry et al.(2013)] BAKRY, D., GENTIL, I. and LEDOUX, M. (2013). *Analysis and Geometry of Markov*  
909 *Diffusion Operators*. Springer.
- 910 [Bobkov(1999)] BOBKOV, S. G. (1999). Isoperimetric and analytic inequalities for log-concave probability  
911 measures. *Ann. Probab.* **27** 1903–1921.
- 912 [Cattiaux and Guillin(2009)] CATTIAUX, P. and GUILLIN, A. (2009). Trends to equilibrium in total variation  
913 distance. *Ann. Inst. H. Poincaré Probab. Statist.* **45** 117–145.
- 914 [Diaconis and Stroock(1991)] DIACONIS, P. and STROOCK, D. (1991). Geometric bounds for eigenvalues of  
915 markov chains. *The Annals of Applied Probability* **1** 36–61.

- 916 [Dwivedi et al.(2018)] DWIVEDI, R., CHEN, Y., WAINWRIGHT, M. J. and YU, B. (2018). Log-concave sam-  
 917 pling: Metropolis-Hastings algorithms are fast! *arXiv e-prints* arXiv:1801.02309.
- 918 [Frieze et al.(1994)] FRIEZE, A., KANNAN, R. and POLSON, N. (1994). Sampling from log-concave distribu-  
 919 tions. *Ann. Appl. Probab.* **4** 812–837.
- 920 [Ge et al.(2018)] GE, R., LEE, H. and RISTESKI, A. (2018). Simulated Tempering Langevin Monte Carlo II:  
 921 An Improved Proof using Soft Markov Chain Decomposition. *arXiv e-prints* arXiv:1812.00793.
- 922 [George and McCulloch(1997)] GEORGE, E. I. and McCULLOCH, R. E. (1997). Approaches to bayesian vari-  
 923 able selection. *Statist. Sinica* **7** 339–373.
- 924 [Laurent and Massart(2000)] LAURENT, B. and MASSART, P. (2000). Adaptive estimation of a quadratic  
 925 functional by model selection. *Ann. Statist.* **28** 1302–1338.
- 926 [Liggett(1991)] LIGGETT, T. M. (1991).  $l_2$  rates of convergence for attractive reversible nearest particle sys-  
 927 tems: The critical case. *Ann. Probab.* **19** 935–959.
- 928 [Liu(1994)] LIU, J. S. (1994). The collapsed gibbs sampler in bayesian computations with applications to a  
 929 gene regulation problem. *Journal of the American Statistical Association* **89** 958–966.
- 930 [Lovász(1999)] LOVÁSZ, L. (1999). Hit-and-run mixes fast. *Math. Program.* **86** 443–461.
- 931 [Lovász and Simonovits(1993)] LOVÁSZ, L. and SIMONOVITS, M. (1993). Random walks in a convex body and  
 932 an improved volume algorithm. *Random Structures Algorithms* **4** 359–412.
- 933 [Lovász and Vempala(2007)] LOVÁSZ, L. and VEMPALA, S. (2007). The geometry of logconcave functions and  
 934 sampling algorithms. *Random Structures Algorithms* **30** 307–358.
- 935 [Madras and Randall(2002)] MADRAS, N. and RANDALL, D. (2002). Markov chain decomposition for conver-  
 936 gence rate analysis. *Ann. Appl. Probab.* **12** 581–606.
- 937 [Meinshausen and Yu(2009)] MEINSHAUSEN, N. and YU, B. (2009). Lasso-type recovery of sparse representa-  
 938 tions for high-dimensional data. *Ann. Statist.* **37** 246–270.
- 939 [Montenegro and Tetali(2006)] MONTENEGRO, R. and TETALI, P. (2006). Mathematical aspects of mixing  
 940 times in markov chains. *Found. Trends Theor. Comput. Sci.* **1** 237–354.
- 941 [Narisetty and He(2014)] NARISSETTY, N. and HE, X. (2014). Bayesian variable selection with shrinking and  
 942 diffusing priors. *Ann. Statist.* **42** 789–817.
- 943 [Raginsky et al.(2017)] RAGINSKY, M., RAKHLIN, A. and TELGARSKY, M. (2017). Non-convex learning via  
 944 stochastic gradient langevin dynamics: a nonasymptotic analysis. In *Proceedings of the 2017 Confer-*  
 945 *ence on Learning Theory*, vol. 65 of *Proceedings of Machine Learning Research*. PMLR.
- 946 [Raskutti et al.(2010)] RASKUTTI, G., WAINWRIGHT, M. J. and YU, B. (2010). Restricted eigenvalue proper-  
 947 ties for correlated gaussian designs. *J. Mach. Learn. Res.* **11** 2241–2259.
- 948 [Sinclair(1992)] SINCLAIR, A. (1992). Improved bounds for mixing rates of markov chains and multicommodity  
 949 flow. *Combinatorics, Probability and Computing* **1** 351–370.
- 950 [Welling and Teh(2011)] WELLING, M. and TEH, Y. W. (2011). Bayesian learning via stochastic gradient  
 951 langevin dynamics. In *Proceedings of the 28th International Conference on International Conference*  
 952 *on Machine Learning*. ICML’11, Omnipress, USA.  
 953 URL <http://dl.acm.org/citation.cfm?id=3104482.3104568>
- 954 [Yang et al.(2016)] YANG, Y., WAINWRIGHT, M. J. and JORDAN, M. I. (2016). On the computational com-  
 955 plexity of high-dimensional bayesian variable selection. *Ann. Statist.* **44** 2497–2532.