

from $p(\mathbf{x})$ would be completely independent from each other and would visit many different regions in \mathbf{x} space proportional to their probability. Instead, especially in high dimensional cases, MCMC samples become very correlated. We refer to such behavior as slow mixing or even failure to mix. MCMC methods with slow mixing can be seen as inadvertently performing something resembling noisy gradient descent on the energy function, or equivalently noisy hill climbing on the probability, with respect to the state of the chain (the random variables being sampled). The chain tends to take small steps (in the space of the state of the Markov chain), from a configuration $\mathbf{x}^{(t-1)}$ to a configuration $\mathbf{x}^{(t)}$, with the energy $E(\mathbf{x}^{(t)})$ generally lower or approximately equal to the energy $E(\mathbf{x}^{(t-1)})$, with a preference for moves that yield lower energy configurations. When starting from a rather improbable configuration (higher energy than the typical ones from $p(\mathbf{x})$), the chain tends to gradually reduce the energy of the state and only occasionally move to another mode. Once the chain has found a region of low energy (for example, if the variables are pixels in an image, a region of low energy might be a connected manifold of images of the same object), which we call a mode, the chain will tend to walk around that mode (following a kind of random walk). Once in a while it will step out of that mode and generally return to it or (if it finds an escape route) move towards another mode. The problem is that successful escape routes are rare for many interesting distributions, so the Markov chain will continue to sample the same mode longer than it should.

This is very clear when we consider the Gibbs sampling algorithm (section 17.4). In this context, consider the probability of going from one mode to a nearby mode within a given number of steps. What will determine that probability is the shape of the “energy barrier” between these modes. Transitions between two modes that are separated by a high energy barrier (a region of low probability) are exponentially less likely (in terms of the height of the energy barrier). This is illustrated in figure 17.1. The problem arises when there are multiple modes with high probability that are separated by regions of low probability, especially when each Gibbs sampling step must update only a small subset of variables whose values are largely determined by the other variables.

As a simple example, consider an energy-based model over two variables a and b , which are both binary with a sign, taking on values -1 and 1 . If $E(a, b) = -wab$ for some large positive number w , then the model expresses a strong belief that a and b have the same sign. Consider updating b using a Gibbs sampling step with $a = 1$. The conditional distribution over b is given by $P(b = 1 \mid a = 1) = \sigma(w)$. If w is large, the sigmoid saturates, and the probability of also assigning b to be 1 is close to 1 . Likewise, if $a = -1$, the probability of assigning b to be -1 is close to 1 . According to $P_{\text{model}}(a, b)$, both signs of both variables are equally likely.