# Markov chain Monte Carlo methods

- The approximations of expectation so far have assumed that the samples are independent draws.

- This sounds good, but in high dimensions, we do not know how to get **good** independent samples from the distribution.

- MCMC methods drop this requirement.

- Basic intuition
  – If you have **finally** found a region of high probability, stick around for a bit, enjoy yourself, grab some more samples.

# Markov chain Monte Carlo methods

- Samples are conditioned on the previous one (this is the Markov chain).

- MCMC is generally a good hammer for complex, high dimensional, problems.

- Main downside is that it is not "plug-and-play"
  – Doing well requires taking advantage to the structure of your problem

  – MCMC tends to be expensive (but take heart---there may not be any other solution, and at least your problem is being solved).

# Metropolis Example

We want samples $z^{(1)}, z^{(2)}, \ldots$

Again, write $p(z) = \tilde{p}(z)/Z$

Assume that $q\left(z \middle| z^{(prev)}\right)$ can be sampled easily

Also assume that $q(\ )$ is symmetric, i.e., $q\left(z_A \middle| z_B\right) = q\left(z_B \middle| z_A\right)$

For example, $q\left(z \middle| z^{(prev)}\right) \sim \mathbb{N}\left(z; z^{(prev)}, \sigma^2\right)$

# Metropolis Example

While not_bored
{

    Sample $q\left(z \middle| z^{(prev)}\right)$

    Accept with probability $A\left(z, z^{(prev)}\right) = \min\left(1, \dfrac{\tilde{p}(z)}{\tilde{p}\left(z^{(prev)}\right)}\right)$

    If accept, emit $z$, otherwise, emit $z^{(prev)}$.

}

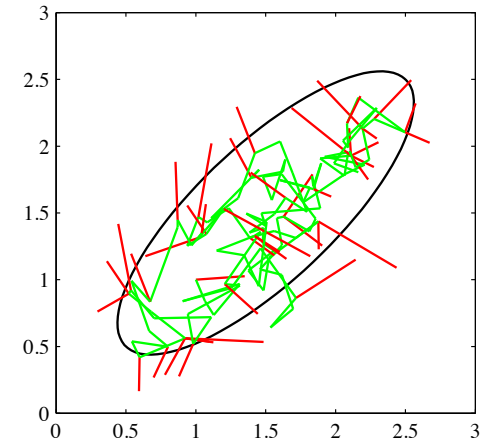If things get better, always accept. If they get worse, sometimes accept.

## Metropolis Example

Note that

$$A\left(z, z^{(prev)}\right) = \min\left(1, \frac{\tilde{p}(z)}{\tilde{p}\left(z^{(prev)}\right)}\right) = \min\left(1, \frac{p(z)}{p\left(z^{(prev)}\right)}\right)$$

We do not need to normalize $p(z)$

## Metropolis Example



Green follows accepted proposals
Red are rejected moves.

## Markov chain view

Denote an initial probability distribution by $p\left(z^{(1)}\right)$

Define transition probabilities by:

$$T\left(z^{(prev)}, z\right) = p\left(z \mid z^{(prev)}\right) \qquad \text{(a probability distribution)}$$

$T = T_m(\ )$ can change over time, but for now, assume that it
it is always the same (homogeneous chain)

A given chain evolves from a sample of $p\left(z^{(1)}\right)$, and is
an instance from an essemble of chains.

## Stationary Markov chains

- Recall that our goal is to have our Markov chain emit samples from our target distribution.

- This implies that the distribution being sampled at time $t+1$ would be the same as that of time $t$ (stationary).

- If our stationary (target) distribution is $p()$, then if imagine an ensemble of chains, they are in each state with (long-run) probability $p()$.
  - On average, a switch from s1 to s2 happens as often as going from s2 to s1, otherwise, the percentage of states would not be stable

- If our stationary (target) distribution is $p()$, what do the transition probabilities look like?

# Detailed balance

- Detailed balance is defined by:

$$p(z)T(z,z') = p(z')T(z',z)$$

  (We assume that $T(\cdot) > 0$)

- Detailed balance is a sufficient condition for a stationary distribution.

- Detailed balance is also referred to as reversibility.

# Detailed balance implies stationary

$$p(z) = \sum_{z'} p(z')T(z',z) \qquad \text{(marginalization)}$$

$$p(z')T(z',z) = p^{(prev)}(z)T(z,z') \qquad \text{(assuming detailed balance)}$$

$$p(z) = \sum_{z'} p(z')T(z',z) = \sum_{z'} p^{(prev)}(z)T(z,z') = p^{(prev)}(z)\underbrace{\sum_{z'} T(z,z')}_{\text{This is 1}} = p^{(prev)}(z)$$

Pedantically, $\qquad \underbrace{\sum_{z'} T(z,z') = \sum_{z'} p(z'|z) = \sum_{z'} \dfrac{p(z',z)}{p(z)} = \dfrac{p(z)}{p(z)} = 1}_{\text{Always true (a conditional probability is a probability)}}$

Hence, detailed balance implies the distribution is stationary.

# Detailed balance (cont)

- Detailed balance (for $p()$) means that *if* our chain was generating samples from $p()$, it would continue to due so.
  - We will address how it gets there shortly

- Does the Metropolis algorithm have detailed balance?

# Metropolis Example

While not_bored
{

    Sample $q\left(z \middle| z^{(prev)}\right)$

    Accept with probability $A\left(z, z^{(prev)}\right) = \min\left(1, \dfrac{\tilde{p}(z)}{\tilde{p}\left(z^{(prev)}\right)}\right)$

    If accept, emit $z$, otherwise, emit $z^{(prev)}$.

}

Same as $\dfrac{p(z)}{p\left(z^{(prev)}\right)}$

## Metropolis Example

Recall that in Metropolis, $\quad A(z,z') = \min\left(1, \dfrac{p(z)}{p(z')}\right)$

For detailed balance, we need to show

$$p(z')q(z|z')A(z,z') = p(z)q(z'|z)A(z',z)$$

> Probability of transition from $z'$ to $z$ is the probability that $z'$ is proposed, **and** it is accepted.

## Metropolis Example

Recall that in Metropolis, $\quad A(z,z') = \min\left(1, \dfrac{p(z)}{p(z')}\right)$

$$p(z')q(z|z')A(z,z') = q(z|z')\min\big(p(z'),p(z)\big)$$
$$= q(z'|z)\min\big(p(z'),p(z)\big)$$
$$= p(z)q(z'|z)\min\left(\frac{p(z')}{p(z)},1\right)$$
$$= p(z)q(z'|z)\min\left(1,\frac{p(z')}{p(z)}\right)$$
$$= p(z)q(z'|z)A(z',z)$$

## Ergodic chains

- Different starting probabilities will give different chains

- We want our chains to converge (in the limit) to the same stationary state, regardless of starting distribution.

- Such chains are called ergodic, and the common stationary state is called the equilibrium state.

- Ergodic chains have a unique equilibrium.

## When do our chains converge?

- Important theorem tells us that (for finite state spaces*) our chains converge to equilibrium under two relatively weak conditions.

- (1) Irreducible
  - We can get from any state to any other state
- (2) Aperiodic
  - The chain does not get trapped in cycles

- These are true for detailed balance with T>0 which is sufficient, but not necessary for convergence.

*Infinite or uncountable state spaces introduces additional complexities.

## Evolution of ergodic chains

Let $p^{(t)}(z)$ be the distribution at some time (e.g., initial distribution)

Let $\pi(z)$ be the stationary distribution

Let $p^{(t)}(z) = \pi(z) - \Delta^{(t)}(z)$

Note that the elements of $p^{(t+1)}(z)$ and $\pi(z)$ sum to one, and thus the elements of $\Delta(z)$ sum to zero.

Note also that $\Delta(z)$ is not a probablity.

## Evolution of ergodic chains

Let $p^{(t)}(z)$ be the distribution at some time (e.g., initial distribution)

Let $\pi(z)$ be the stationary distribution

Let $p^{(t)}(z) = \pi(z) - \Delta^{(t)}(z)$

$$p^{(t+1)}(z) = \sum_{z'} p^{(t)}(z') \, T(z,z')$$
$$= \sum_{z'} \pi(z') \, T(z,z') - \sum_{z'} \Delta^{(t)}(z') \, T(z,z')$$
$$= \pi(z) - \Delta^{(t+1)}(z)$$

## Evolution of ergodic chains

$$p^{(t+1)}(z) = \sum_{z'} p^{(t)}(z') \, T(z,z')$$
$$= \sum_{z'} \pi(z') \, T(z,z') - \sum_{z'} \Delta^{(t)}(z') \, T(z,z')$$
$$= \pi(z) - \Delta^{(t+1)}(z)$$

Claim that $\left| \Delta^{(t)}(z) \right| < (1-v)^t$

where $\quad v = \min_{z} \min_{z' : \pi(z') > 0} \dfrac{T(z,z')}{\pi(z)}$

and we have $\quad 0 < v \leq 1$

(see final, optional problem #5)