

# A Haskell Perspective on a general perspective on the Metropolis–Hastings kernel

Dominic Steinitz

November 14, 2022

## Todo list

I think I have squared some of this but need to write this down . . . . .	8
This needs a bit more explanation . . . . .	9
This needs a bit more explanation . . . . .	10
It should make no difference what we return in the momentum position; yet it does? Maybe not but at least investigate it. . . . .	11
It’s not clear this is a useful section but who knows? . . . . .	11

## 1 Introduction

Remarkably (to me at least) all<sup>1</sup> MCMC algorithms can be captured by one general algorithm. At the moment you are expected to know how MCMC works to be able to read what follows. I may add a section introducing MCMC later.

Here’s **Algorithm 1** from [1]:

```
algo1 :: Show a => Show b => (MonadDistribution m, Fractional t) =>
  (a, c) -> (a -> m b) -> ((a, b) -> (a, b)) -> (t -> Double) -> ((a, b) -> t) -> m (a, b)
algo1 (ξ0, -) μξ0 φ a ρ = do
  μξ-0 ← μξ0 ξ0
  let ξ = (ξ0, μξ-0)
  let α = a $ (ρ ∘ φ) ξ / ρ ξ
  u ← random
  if u < α
  then return $ φ ξ
  else return ξ
```

Something very similar seems to have been discovered in [3] and [5]. Serendipitously, all three papers call this algorithm 1!

---

<sup>1</sup>well almost all

## 2 An Example with an Analytical Solution

In Bayesian statistics we have a prior distribution for the unknown mean which we also take to be normal

$$\mu \sim \mathcal{N}(\mu_0, \sigma_0^2)$$

and then use a sample

$$x \mid \mu \sim \mathcal{N}(\mu, \sigma^2)$$

to produce a posterior distribution for it

$$\mu \mid x \sim \mathcal{N}\left(\frac{\sigma_0^2}{\sigma^2 + \sigma_0^2}x + \frac{\sigma^2}{\sigma^2 + \sigma_0^2}\mu_0, \left(\frac{1}{\sigma_0^2} + \frac{1}{\sigma^2}\right)^{-1}\right)$$

If we continue to take samples then the posterior distribution becomes

$$\mu \mid x_1, x_2, \dots, x_n \sim \mathcal{N}\left(\frac{\sigma_0^2}{\frac{\sigma^2}{n} + \sigma_0^2}\bar{x} + \frac{\sigma^2}{\frac{\sigma^2}{n} + \sigma_0^2}\mu_0, \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)^{-1}\right)$$

Note that if we take  $\sigma_0$  to be very large (we have little prior information about the value of  $\mu$ ) then

$$\mu \mid x_1, x_2, \dots, x_n \sim \mathcal{N}\left(\bar{x}, \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}\right)^{-1}\right)$$

and if we take  $n$  to be very large then

$$\mu \mid x_1, x_2, \dots, x_n \sim \mathcal{N}\left(\bar{x}, \frac{\sigma}{\sqrt{n}}\right)$$

which ties up with the classical estimate.

Let's illustrate this with a few numbers.

```
 $\mu_0, \sigma_0, \sigma, \sigma_P, z :: \text{Floating } a \Rightarrow a$   
 $\mu_0 = 0.0$   
 $\sigma_0 = 1.0$   
 $\sigma = 1.0$   
 $\sigma_P = 0.2$   
 $z = 4.0$   
 $\hat{\mu} :: \text{Double}$   
 $\hat{\mu} = z * \sigma_0 \uparrow 2 / (\sigma \uparrow 2 + \sigma_0 \uparrow 2) + \mu_0 * \sigma \uparrow 2 / (\sigma \uparrow 2 + \sigma_0 \uparrow 2)$   
 $\hat{\sigma} :: \text{Double}$   
 $\hat{\sigma} = \text{sqrt } \$ \text{ recip } (\text{recip } \sigma_0 \uparrow 2 + \text{recip } \sigma \uparrow 2)$ 
```

This gives  $\hat{\mu} = 2.0$  and  $\hat{\sigma} = 0.7071067811865476$  which is what we would expect: we thought the mean was  $\mu_0 = 0.0$  but we have an observation  $z = 4.0$  and also the variance is now less.

### 3 Using MCMC

For us, we want the posterior

$$\varpi(\mu) = \frac{1}{Z} \exp \frac{(x - \mu)^2}{2\sigma^2} \exp \frac{(\mu - \mu_0)^2}{2\sigma_0^2}$$

where  $x, \mu_0, \sigma$  and  $\sigma_0$  are all given but  $Z$  is unknown.

#### 3.1 Random Walk Metropolis

Let's implement a traditional random walk. Here's the proposal distribution:

```
Q :: Double -> Double -> Double
Q w w' = exp (-(w - w') ↑ 2 / (2 * σ_P ↑ 2))
```

And here's the specification for  $\rho$ :

```
ρ̃ :: (a -> Double) -> (a -> b -> Double) -> (a, b) -> Double
ρ̃ φ q (w, w') = φ w * q w w'
```

Here's the un-normalised posterior:

```
φ̃ :: Floating a => a -> a
φ̃ μ = exp (-(z - μ) ↑ 2 / (2 * σ ↑ 2)) * exp (-(μ - μ_0) ↑ 2 / (2 * σ_0 ↑ 2))
```

We can now use one step of the algorithm and then run it for as many times as we wish:

```
testRwmOneStep :: MonadDistribution m => (Double, Double) -> m (Double, Double)
testRwmOneStep (ξ_0, -) = algo1 (ξ_0, ⊥) μ_ξ_0 φ a ρ
  where
    φ = λ(x, y) -> (y, x)
    a = min 1.0
    ρ = ρ̃ φ Q
    μ_ξ_0 = λζ -> normal ζ σ_P
testRwm :: (Eq a, Num a, MonadDistribution m) =>
  a -> m [(Double, Double)]
testRwm n = unfoldM f (n, (1.0, 0.0 / 0.0))
  where
    f (0, -) = return Nothing
    f (m, s) = do x ← testRwmOneStep s
                  return $ Just (s, (m - 1, x))
```

And we can see the results in Figure 1. A bit skewed but we didn't burn in and the starting value is 1.0.

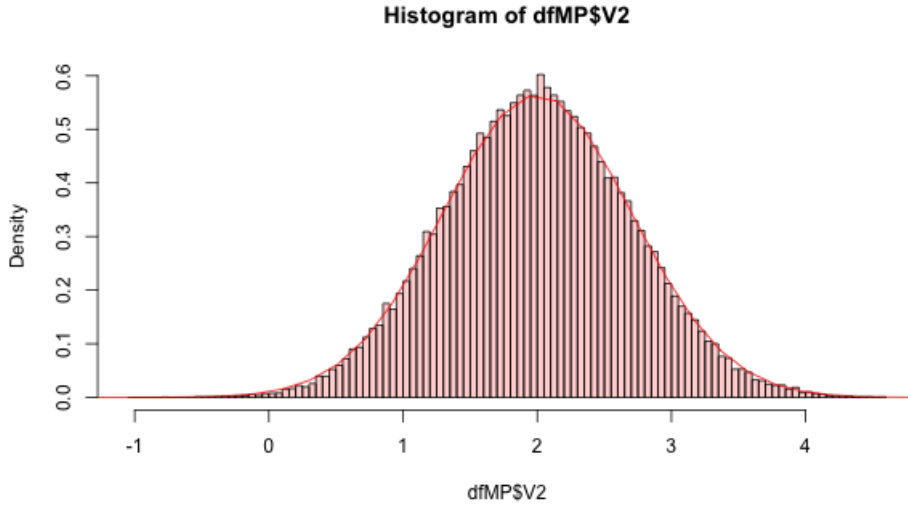


Figure 1: Random Walk Metropolis

### 3.2 Random walk Metropolis ratio

Here's a different algorithm expressed using the generalised approach. The results are in Figure 2.

```
testMwMrOneStep :: MonadDistribution m => (Double, Double) -> m (Double, Double)
```

```
testMwMrOneStep (ξ0, ⊥) = algo1 (ξ0, ⊥) μξ0 φ a ρ
```

**where**

```
φ = λ(x, y) -> (x + y, -y)
```

```
a = min 1.0
```

```
ρ =  $\tilde{\rho}$   $\tilde{\varphi}$  (\_ -> \_ -> 1.0)
```

```
μξ0 = const (quantile (normalDistr 0.0 1.0) < $ > random)
```

```
testMwMr :: (Eq a, Num a, MonadDistribution m) =>
```

```
a -> m [(Double, Double)]
```

```
testMwMr n = unfoldM f (n, (1.0, 0.0 / 0.0))
```

**where**

```
f (0, _) = return Nothing
```

```
f (m, s) = do x ← testMwMrOneStep s
```

```
return $ Just (s, (m - 1, x))
```

### 3.3 What monad-bayes does

Here's our toy problem expressed in monad-bayes:

```
singleObs :: (MonadDistribution m, MonadFactor m) => m Double
```

```
singleObs = do
```

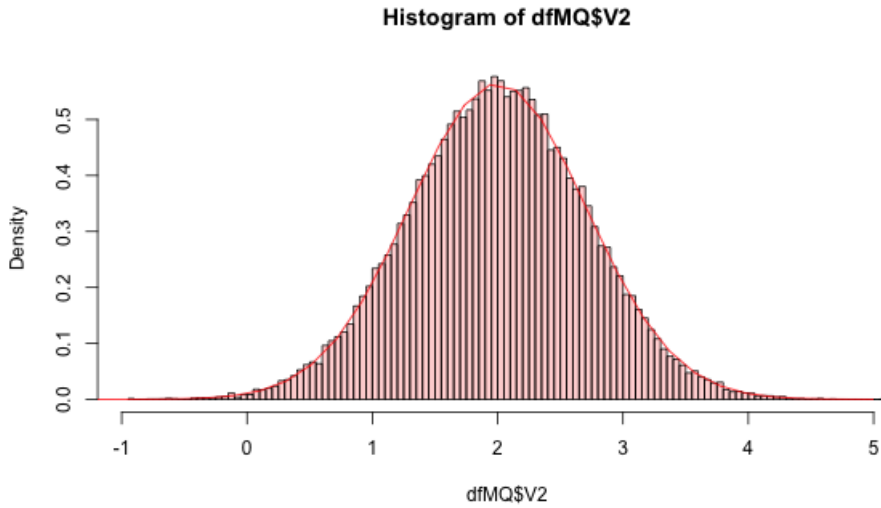


Figure 2: Random walk Metropolis ratio

```

μ ← normal μ₀ σ₀
factor $ normalPdf μ σ z
return μ

```

Here's what I think monad-bayes does with this using the General Perspective. The results are in Figure 3.

```

testMbOneStep :: MonadDistribution m ⇒ (Double, Double) → m (Double, Double)
testMbOneStep (ξ₀, -) = algo1 (ξ₀, ⊥) μξ₀ φ a ρ
  where
    φ = λ(x, y) → (y, x)
    a = min 1.0
    ρ =  $\tilde{\rho}$  (λμ → exp (-(z - μ) ↑ 2 / (2 * σ ↑ 2))) (\_ → \_ → 1.0)
    μξ₀ = const (quantile (normalDistr μ₀ σ₀) < $ > random)
testMb :: (Eq a, Num a, MonadDistribution m) ⇒
  a → m [(Double, Double)]
testMb n = unfoldM f (n, (1.0, 0.0 / 0.0))
  where
    f (0, -) = return Nothing
    f (m, s) = do x ← testMbOneStep s
      return $ Just (s, (m - 1, x))

```

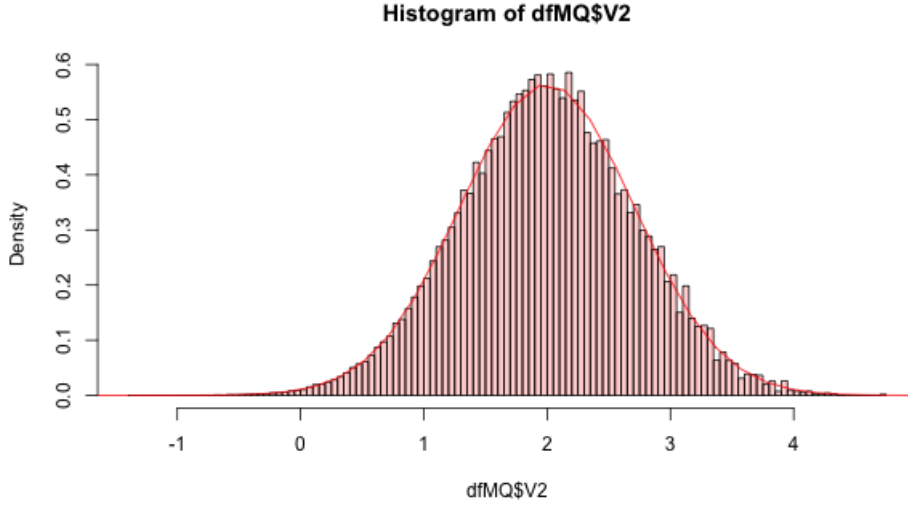


Figure 3: monad-bayes

## 4 Some Mathematical Notes

Suppose we don't know the classical MCMC algorithm. We can derive it from [1]:

$$r(z, z') = \frac{\varpi(z') q(z', z)}{\varpi(z) q(z, z')}$$

But where does this come from? We define  $\mu$ :

$$\mu(d\xi) \triangleq \pi(d\xi_0) \mu_{\xi_0}(d\xi_{-0})$$

$$\mu(d(z, z')) \triangleq \varpi(z) dz q_z(z') dz'$$

Let  $\mu$  be a finite measure on  $(E, \mathcal{E})$ ,  $\phi: E \rightarrow E$  an involution, let  $\lambda \gg \mu$  be a  $\sigma$ -finite measure satisfying  $\lambda \equiv \lambda^\phi$  and let  $\rho = d\mu/d\lambda$ . Then we can take  $S = S(\mu, \mu^\phi)$  to be  $S = \{\xi : \rho(\xi) \wedge \rho \circ \phi(\xi) > 0\}$  and

$$r(\xi) = \begin{cases} \frac{\rho \circ \phi(\xi) \frac{d\lambda^\phi}{d\lambda}(\xi)}{\rho(\xi)} & \xi \in S, \\ 0 & \text{otherwise} \end{cases}$$

So

$$\rho(z, z') \triangleq \varpi(z) q_z(z')$$

and with  $\phi(z, z') = (z', z)$  we regain the familiar

$$r(z, z') = \frac{\varpi(z')q(z', z)}{\varpi(z)q(z, z')}$$

## 5 Student's T

Let's try running it on Student's T with 5 degrees of freedom using what I hope the textbook presentation of Metropolis—Hastings. The probability density function (aka the Radon-Nikodym derivative wrt Lebesgue measure) is

$$f(t) = \frac{8}{3\pi\sqrt{5} \left(1 + \frac{t^2}{5}\right)^3}$$

It's traditional to have  $q_z(\cdot) \sim \mathcal{N}(z, \sigma_p^2)$  for some given  $\sigma_p$ .

Here's the density function for Student's T with 5 degrees of freedom. We've defined it in terms of an un-normalised density so that we can pretend we don't know the normalisation constant but still sample from the distribution via MCMC.

```
student5U :: Floating a => a -> a
student5U t = 1 / (1 + t ↑ 2 / 5) ↑ 3
student5 :: Floating a => a -> a
student5 t = student5U t * 8 / (3 * pi * sqrt 5)
```

Again, we can now use one step of the algorithm and then run it for as many times as we wish. We instantiate the algorithm to be a Random Walk Metropolis.

```
testStudentRwmOneStep :: MonadDistribution m =>
  (Double, Double) -> m (Double, Double)
testStudentRwmOneStep (ξ₀, -) = algo1 (ξ₀, ⊥) (λζ -> normal ζ σP)
  (λ(x, y) -> (y, x)) (min 1.0) (ρ student5U Q)
testStudentRwm :: (Eq a, Num a, MonadDistribution m) =>
  a -> m [(Double, Double)]
testStudentRwm n = unfoldM f (n, (0.0, 0.0 / 0.0))
  where
    f (0, -) = return Nothing
    f (m, s) = do x ← testStudentRwmOneStep s
      return $ Just (s, (m - 1, x))
```

We can also instantiate it with what I think monad-bayes does.

```
testStudentMbOneStep :: MonadDistribution m => (Double, b) -> m (Double, Double)
testStudentMbOneStep (ξ₀, -) = algo1 (ξ₀, ⊥) (const ((quantile (studentT 5)) < $ > random))
```

```

( $\lambda(x, y) \rightarrow (x + y, -y)$ ) (min 1.0) (student5U  $\circ$  fst)
testStudentMb :: (Eq a, Num a, MonadDistribution m) =>
  a -> m [(Double, Double)]
testStudentMb n = unfoldM f (n, (0.0, 0.0 / 0.0))
where
  f (0, _) = return Nothing
  f (m, s) = do x <- testStudentMbOneStep s
    return $ Just (s, (m - 1, x))

```

The results are shown in Figure 4 and Figure 5.

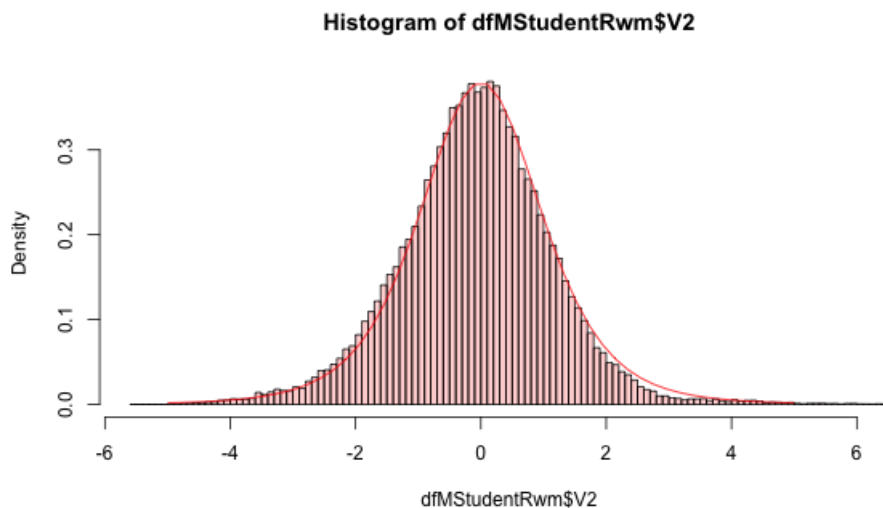


Figure 4: Random Walk Metropolis Student's T 5

## 6 Hamiltonian Monte Carlo

We'd like to put HMC into the same general framework but at the moment, I am having trouble squaring Example 14 in [1] with the algorithm given in [4] (and I haven't even looked in [2]). Here's as far as I got with Student's t-distribution of degree 5. There's something going on with exponentiating the Hamiltonian which I don't understand yet either.

Here's Student's T again:

$$f(t) = \frac{8}{3\pi\sqrt{5} \left(1 + \frac{t^2}{5}\right)^3}$$

I think I have squared some of this but need to write this down



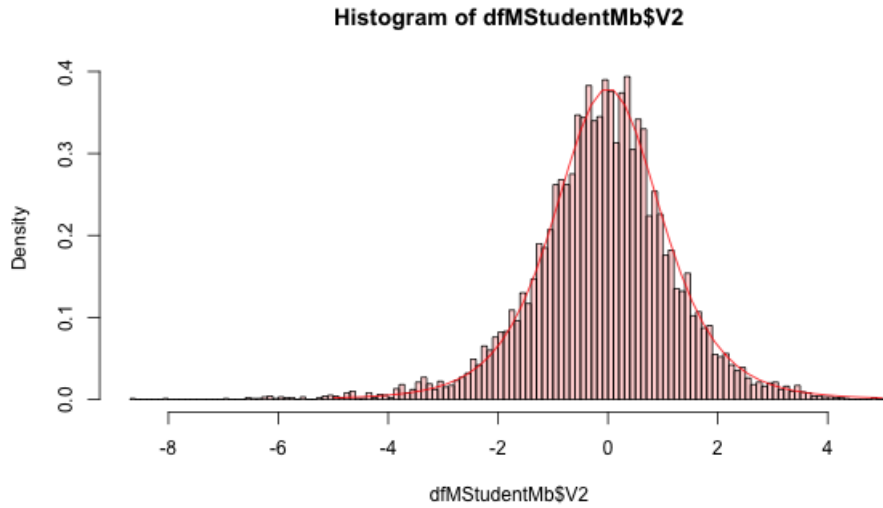


Figure 5: Monad Bayes Student's T 5

Unnormalised:

$$g(t) = \frac{1}{\left(1 + \frac{t^2}{5}\right)^3}$$

And as the potential energy part of the Hamiltonian:

$$U(t) = -\log g(t) = 3 \log \left(1 + \frac{r^2}{5}\right)$$

Here's a version of the leapfrog algorithm:

```

leapfrog :: Fractional a => a -> Int -> (a -> a) -> (a, a) -> (a, a)
leapfrog epsilon l gradU (qPrev, p) = (q1, p3)
  where
    p' = p - epsilon * gradU qPrev / 2
    f 0 (qOld, pOld) = r
      where
        qNew = qOld + epsilon * pOld
        pNew = pOld
        r = (qNew, pNew)
    f _ (qOld, pOld) = r
      where

```

This needs a bit more explanation

```

qNew = qOld + epsilon * pOld
pNew = pOld - epsilon * gradU qNew
r = (qNew, pNew)
(q1, p1) = foldr f (qPrev, p') ([0..l-1])
p2 = p1 - epsilon * gradU q1 / 2
-- Is this necessary?
p3 = negate p2

```

This is the Hamiltonian:

```

rhoHmc :: Floating a => (b -> a) -> (b, a) -> a
rhoHmc u (q, p) = pU * pK
  where
    pU = recip $ u q
    pK = exp $ p ↑ 2 / 2

```

This needs a bit more explanation

We need the derivative of the potential energy for the leapfrog method. We could use automatic differentiation of course.

```

gradU :: Fractional a => a -> a
gradU r = 3 * (2 * r / 5) / (1 + (r ↑ 2) / 5)
bigU :: Floating a => a -> a
bigU = negate ∘ log ∘ student5U
gradUAD :: Floating a => a -> a
gradUAD w = case grad (λ[x] -> bigU x) $ [w] of
  [y] -> y
  _ -> error "Whatever"

```

And now we can run the sampler. The results are in 6.

```

eta :: Fractional a => a
eta = 0.3
bigL :: Int
bigL = 10
testHmcOneStep :: MonadDistribution m => (Double, Double) -> m (Double, Double)
testHmcOneStep (ξ0, -) = algo1 (ξ0, ⊥) μξ0 φ a ρ
  where
    φ = leapfrog eta bigL gradU
    a = min 1.0
    ρ = rhoHmc student5U
    μξ0 = const $ normal 0.0 1.0
testHmc :: (Eq a, MonadDistribution m, Num a) =>
  a -> m [(Double, Double)]
testHmc n = unfoldM f (n, (0.0, 0.0))
  where

```

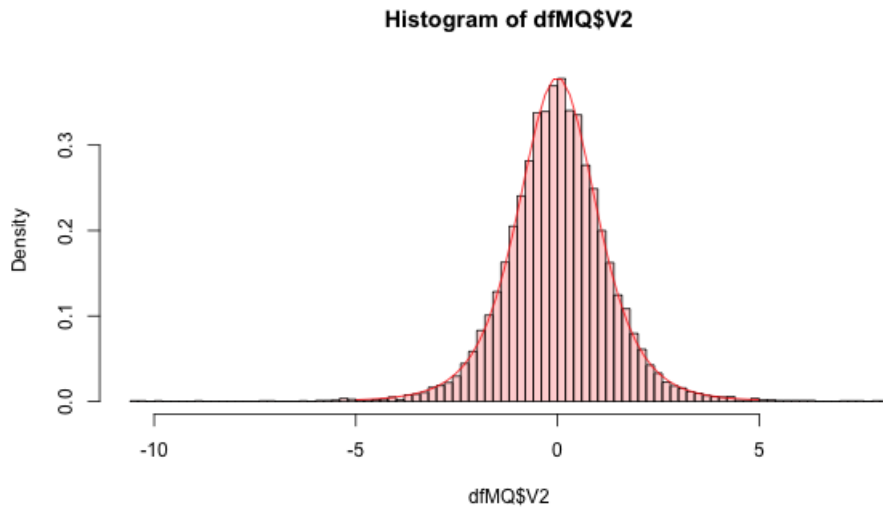


Figure 6: Hamiltonian Monte Carlo Student's T 5

```
f (0, _) = return Nothing
f (m, s) = do a ← testHmcOneStep s
  return $ Just (s, (m - 1, a))
```

## 7 Gen

Gen is a probabilistic programming language. I've taken the example from [3] and converted it to use monad-bayes.

```
genEg :: MonadDistribution m => Int -> m [Double]
genEg n = do
  k ← (+1) < $ > poisson 1.0
  means ← replicate k < $ > normal 0.0 10.0
  gammas ← replicate k < $ > gamma 1.0 10.0
  let invGammas = map recip gammas
  weights ← dirichlet (V.replicate k 2.0)
  replicate n < $ > (categorical weights >>= λi → normal (means !! i) (invGammas !! i))
```

It should make no difference what we return in the momentum position; yet it does? Maybe not but at least investigate it.

It's not clear this is a useful section but who knows?

## 8 Bibliography

## References

- [1] Christophe Andrieu, Anthony Lee, and Sam Livingstone. A general perspective on the metropolis-hastings kernel. *arXiv*, December 2020. [1](#), [6](#), [8](#)
- [2] M. J. Betancourt, Simon Byrne, Samuel Livingstone, and Mark Girolami. The geometric foundations of hamiltonian monte carlo, 2014. [8](#)
- [3] Marco Cusumano-Towner, Alexander K Lew, and Vikash K Mansinghka. Automating involutive mcmc using probabilistic and differentiable programming. *arXiv preprint arXiv:2007.09871*, 2020. [1](#), [11](#)
- [4] Radford M. Neal. Mcmc using hamiltonian dynamics. *arXiv: Computation*, pages 139–188, 2011. [8](#)
- [5] Kirill Neklyudov, Max Welling, Evgenii Egorov, and Dmitry Vetrov. Involutive mcmc: a unifying framework. In *International Conference on Machine Learning*, pages 7273–7282. PMLR, 2020. [1](#)