

# Effective random generation of algebraic data types through Boltzmann samplers

Maciej Bendkowski

April 4, 2024

# Motivating example: General purpose codecs

Suppose that we want to develop a general purpose codec:

```
1 | class Codec a where
2 |     encode :: a -> ByteString
3 |     decode :: ByteString -> a
```

such that

- ▶ `encode . decode = id`
- ▶ `decode . encode = id`
- ▶ `size (encode a) ≤ size a`

where

```
1 | class Sized a where
2 |     size :: a -> Int
```

# Motivating example: General purpose codecs

Q: How to en- and decode `BinTrees`?

```
1  data BinTree
2      = Node BinTree BinTree
3      | Leaf
4
5  -- / >>> code (Node Leaf (Node Leaf Leaf)) == "10100"
6  code :: BinTree -> String
7  code = \case
8      Leaf -> "0"
9      Node l r -> "1" <> code l <> code r
10
11 uncode :: String -> BinTree
12 uncode = (...)
```

# Motivating example: General purpose codecs

```
1  instance Codec BinTree where
2      encode :: BinTree -> ByteString
3      encode = Binary.encode . LempelZiv.encode . code
4
5      decode :: ByteString -> BinTree
6      decode = uncode . LempelZiv.decode . Binary.decode
7
8  instance Size BinTree where
9      size :: BinTree -> Int
10     size = \case
11         Leaf -> 1
12         Node lt rt -> 1 + size lt + size rt
```

# Motivating example: General purpose codecs

In order to test our implementation, we resort to QuickCheck and generate random **BinTrees** to test the specified invariants:

```
1  newtype ArbBinTree = ArbBinTree BinTree
2  instance Arbitrary ArbBinTree where
3      arbitrary = ArbBinTree <$> genBinTree
4  where
5      genBinTree =
6          frequency
7              [ (1, pure Leaf)
8                , (1, Node <$> genBinTree <*> genBinTree)
9              ]
```

# Motivating example: General purpose codecs

```
1 testProperty "Codec can de- and encode" $
2     \ (ArbBinTree bt) ->
3         Codec.decode (Codec.encode bt) == bt
4
5 testProperty "Encoded trees are compressed" $
6     \ (ArbBinTree bt) ->
7         size (encode s) <= size s
```

Our implementation works, but *how good is it?*

# Motivating example: General purpose codecs

In order to *evaluate* the performance of our prototype codec we can take the average compression ratio over a large corpus of *random* binary trees.

## Problem

- ▶ What **BinTree** distribution should we choose?
- ▶ How to generate **BinTrees** following the chosen distribution?
- ▶ How to deal with more involved algebraic data types?

## Example: Random $\lambda$ -terms





## Example: Random $\lambda$ -terms

```
1  data DeBruijn
2    = Z | S DeBruijn
3
4  data Lambda
5    = Index DeBruijn | Abs Lambda | App Lambda Lambda
6
7  instance Arbitrary DeBruijn where
8    arbitrary =
9      frequency
10        [ (?, pure Z)
11          , (?, S <$> arbitrary)
12        ]
13
14  instance Arbitrary Lambda where
15    arbitrary =
16      frequency
17        [ (?, Index <$> arbitrary)
18          , (?, Abs <$> arbitrary)
19          , (?, App <$> arbitrary <*> arbitrary)
20        ]
```

## Example: Random $\lambda$ -terms

```
1  data DeBruijn
2    = Z | S DeBruijn
3
4  data Lambda
5    = Index DeBruijn | Abs Lambda | App Lambda Lambda
6
7  instance Arbitrary DeBruijn where
8    arbitrary =
9      frequency
10        [ (0.509308127, pure Z)
11          , (0.49069187299999995, S <$> arbitrary)
12        ]
13
14  instance Arbitrary Lambda where
15    arbitrary =
16      frequency
17        [ (0.3703026, Index <$> arbitrary)
18          , (0.25939476, Abs <$> arbitrary)
19          , (0.3703026, App <$> arbitrary <*> arbitrary)
20        ]
```

# Boltzmann samplers

```
1 class BoltzmannSampler a where
2   sample ::
3     RandomGen g
4     => UpperBound
5     -> MaybeT (BitOracle g) (a, Int)
6
7 mkDefBoltzmannSampler ''Lambda 10_000
```

Key properties:

- ▶ the outcome *size* of `sample` is a random variables,
- ▶ instances of `BoltzmannSampler a` with the *exact same size* have the *exact same probability* of being generated,
- ▶ the expected size of the generated instances of `BoltzmannSampler a` follows the user-declared value, such as 10,000 for `Lambda`.

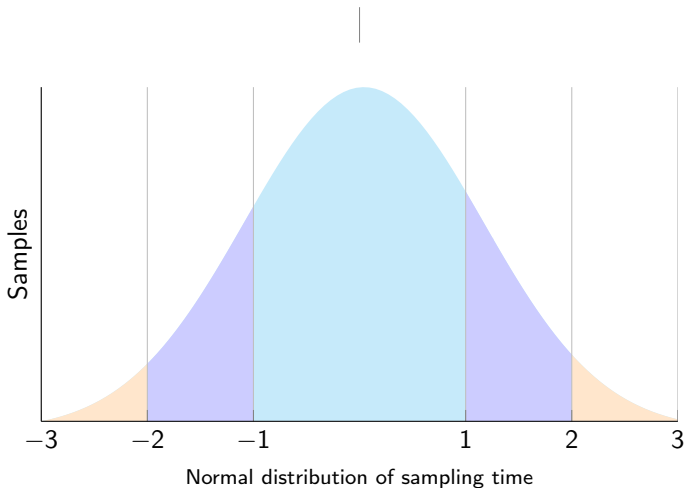
# Boltzmann samplers

With `BoltzmannSampler` a we can have a finer control over the size of outcome samples, e.g.:

```
1 rejectionSampler ::
2   (RandomGen g, BoltzmannSampler a) =>
3     LowerBound -> UpperBound -> BitOracle g a
4 rejectionSampler lb ub = do
5   runMaybeT (sample ub)
6   >>= ( \case
7     Just (obj, s) ->
8       if lb <= s && s <= ub
9       then pure obj
10      else rejectionSampler lb ub
11     Nothing -> rejectionSampler lb ub
12   )
13
```

# Benchmarks

$$\mathcal{N}(\mu = 1.127\text{s}, \sigma^2 = (40.73\text{ms})^2)$$



Sampling time for a default Boltzmann sampler  
generating 1,000 random  $\lambda$ -terms of size in between 800 and 1,200.

# Beyond uniform outcome distribution

We don't have to use the *uniform* outcome distribution!

```
1 mkBoltzmannSampler
2   System
3   { targetType = ''Lambda
4     , meanSize = 10_000
5     , frequencies =
6       ('Abs, 4_000) <:> def
7     , weights =
8       ('Index, 0)
9     <:> $(mkDefWeights ''Lambda)
10  }
```

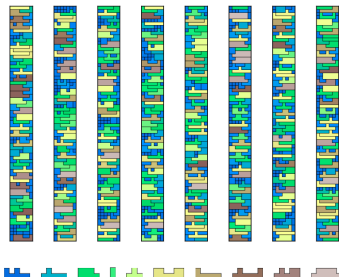


Figure: Five random  $n \times 7$  tilings of areas in the interval  $[500; 520]$  using in total 95 different tiles.

# Generating functions

Let  $\mathcal{S}$  be a set of objects endowed with a *size* function  $|\cdot|: \mathcal{S} \rightarrow \mathbb{N}$  with the property that for all  $n \geq 0$  the set of objects of size  $n$  in  $\mathcal{S}$  is *finite*. For such a class of objects, the corresponding *generating function*  $S(z)$  is a formal power series

$$S(z) = \sum_{n \geq 0} s_n z^n$$

whose coefficients  $(s_n)_{n \geq 0}$  denote the number of objects of size  $n$  in  $\mathcal{S}$ .

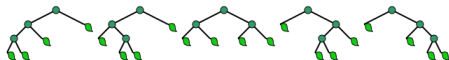
## Example

Fibonacci numbers have a generating function  $F(z)$  of the form

$$F(z) = 1 + z + 2z^2 + 3z^3 + 5z^4 + 8z^5 + 13z^6 + 21z^7 + \dots$$

# Generating functions

```
1 | data BinTree
2 |     = Node BinTree BinTree
3 |     | Leaf
```



The recursive *sum-of-products* nature of algebraic data types lets us find the functional recursion defining the respective generating function:

$$\begin{aligned} C(z) &= 1 + z + 2z^2 + 5z^3 + 14z^4 + 42z^5 + \dots \\ &= 1 + zC(z)^2 \\ &= \frac{1 - \sqrt{1 - 4z}}{2z} \end{aligned}$$

## Remark

Usually, we cannot find closed-form solutions for generating functions.



# Univariate Boltzmann models

Given a real control parameter  $x \in [0, 1]$ , a *Boltzmann model* is a probability distribution in which the probability  $\mathbb{P}_x(\omega)$  of generating an object  $\omega \in \mathcal{S}$  satisfies

$$\mathbb{P}_x(\omega) = \frac{x^{|\omega|}}{S(x)}$$

Note that

- ▶ objects of equal size have equal probabilities, and
- ▶ the outcome size is a varying random variable.

The probability  $\mathbb{P}_x(N = n)$  that the outcome size is  $n$  satisfies

$$\mathbb{P}_x(N = n) = \frac{s_n x^n}{S(x)}$$

whereas

$$\mathbb{E}_x(N) = x \frac{\frac{\partial}{\partial x} S(x)}{S(x)} \quad \sigma_x(N) = \sqrt{x \frac{\frac{\partial}{\partial x} S(x)}{S(x)}}$$

# Parameter tuning

## Problem

Given a system of (possibly mutually recursive) algebraic data types, and a target outcome size  $n$ , find the respective control parameter  $x$  such that

$$n = x \frac{\frac{\partial}{\partial x} S(x)}{S(x)}$$

# Parameter tuning

## Problem

Given a system of (possibly mutually recursive) algebraic data types, and a target outcome size  $n$ , find the respective control parameter  $x$  such that

$$n = x \frac{\frac{\partial}{\partial x} S(x)}{S(x)}$$

## Solution (not covered in this presentation)

Transform the system of partial differential equation into an equivalent *convex optimisation* problem, solve the optimisation problem using efficient interior-point methods, and translate the solution back to the original domain.

# Composition rules for sampler construction

Q: Suppose that we have the control parameter  $x$ .  
How to turn it into an efficient sampler?

- ▶ Let  $\mathcal{S} = \mathcal{A} + \mathcal{B}$  be a co-product class. Then, for  $\omega \in \mathcal{S}$

$$\mathbb{P}_x(\omega \in \mathcal{A}) = \frac{\sum_{\gamma \in \mathcal{A}} x^{|\gamma|}}{S(x)} = \frac{A(x)}{S(x)}$$

Hence, before sampling an object from  $\mathcal{S}$ , we need to make a *skewed coin toss*. With probability  $\frac{A(x)}{S(x)}$  we invoke the sampler for  $\mathcal{A}$ . With probability  $\frac{B(x)}{S(x)}$  we invoke the sampler for  $\mathcal{B}$ .

- ▶ Let  $\mathcal{S} = \mathcal{A} \times \mathcal{B}$  be a product class. Then, for  $\omega = (\alpha, \beta)$

$$\mathbb{P}_x(\omega) = \frac{x^{|\alpha|} x^{|\beta|}}{S(x)} = \frac{x^{|\alpha|}}{A(x)} \frac{x^{|\beta|}}{B(x)} = \mathbb{P}_x(\alpha) \mathbb{P}_x(\beta)$$

Hence, we can invoke both samplers for  $\mathcal{A}$  and  $\mathcal{B}$  and pair up their results.

# Divide and conquer

Compiling Boltzmann samplers breaks down into the following sub-problems:

- ▶ collect the system of ADTs and user parameters (target size, etc.),
- ▶ find all the necessary tuning parameters,
- ▶ compute *branching probabilities* for co-product types,
- ▶ generate *efficient* Haskell code for the sampler.

# Code generation

```
1  instance BoltzmannSampler Lambda where
2    {-# INLINE sample #-}
3    sample !ub =
4      do guard (ub >= 0)
5        lift (BuffonMachine.choice ...)
6          >>=
7          (\case
8            0 -> do (x_0, w_0) <- sample ub
9                  pure (Index x_0, w_0)
10
11           1 -> do (x_0, w_0) <- sample (ub - 1)
12                 (x_1, w_1) <- sample (ub - w_0 - 1)
13                 pure (App x_0 x_1, 1 + w_0 + w_1)
14
15           2 -> do (x_0, w_0) <- sample (ub - 1)
16                 pure (Abs x_0, 1 + w_0)
```

## Sampling from a finite distribution: some ideas

Q: Suppose that we have computed *branching probabilities*

$\mathbb{P} = (p_1, \dots, p_n)$  for a data type with  $n$  distinct constructors. How can we effectively generate from  $\mathbb{P}$ ?

# Sampling from a finite distribution: some ideas

Q: Suppose that we have computed *branching probabilities*

$\mathbb{P} = (p_1, \dots, p_n)$  for a data type with  $n$  distinct constructors. How can we effectively generate from  $\mathbb{P}$ ?

## Solution

Draw a random number  $x \in [0, 1]$  using a standard library call, find the segment  $[p_i, p_{i+1})$  containing  $x$ , and use the respective constructor.



# Sampling from a finite distribution: some ideas

Q: Suppose that we have computed *branching probabilities*  $\mathbb{P} = (p_1, \dots, p_n)$  for a data type with  $n$  distinct constructors. How can we effectively generate from  $\mathbb{P}$ ?

## Solution

Draw a random number  $x \in [0, 1]$  using a standard library call, find the segment  $[p_i, p_{i+1})$  containing  $x$ , and use the respective constructor.

## Better solution

Note that to draw  $x$  we need about 32 random bits whereas we might need much less, e.g. for  $\mathbb{P} = (0.5, 0.5)$ . We can therefore bisect the probability space using one random bit at a time.

# Sampling from a finite distribution: some ideas

Q: Suppose that we have computed *branching probabilities*  $\mathbb{P} = (p_1, \dots, p_n)$  for a data type with  $n$  distinct constructors. How can we effectively generate from  $\mathbb{P}$ ?

## Solution

Draw a random number  $x \in [0, 1]$  using a standard library call, find the segment  $[p_i, p_{i+1})$  containing  $x$ , and use the respective constructor.

## Better solution

Note that to draw  $x$  we need about 32 random bits whereas we might need much less, e.g. for  $\mathbb{P} = (0.5, 0.5)$ . We can therefore bisect the probability space using one random bit at a time.

## Even Better solution

Use the idea of *discrete distribution trees* devised by Knuth and Yao.

# Multiple samplers at a time

```
1  data DeBruijn
2      = Z
3      | S DeBruijn
4
5  data Lambda
6      = Index DeBruijn
7      | Abs Lambda
8      | App Lambda Lambda
9
10 newtype BinLambda = MkBinLambda Lambda
11 mkBoltzmannSampler 'BinLambda 1_000
12
13 -- compiler-generated type
14 newtype Gen_DeBruijn = MkGen_DeBruijn DeBruijn
15
16 -- type coercion(s)
17 coerce
18     @(DeBruijn -> BinLambda)
19     @(Gen_DeBruijn -> BinLambda) Index
```

# Convex optimisation: Random $\lambda$ -terms

```
1 data DeBruijn
2   = Z
3   | S DeBruijn
4
5 data Lambda
6   = Index DeBruijn
7   | Abs Lambda
8   | App Lambda Lambda
```

```
1 import paganini as pg
2 spec = pg.Specification()
3 z = pg.Variable(10000)
4 d = pg.Variable()
5 l = pg.Variable()
6 d = pg.Seq(z)
7
8 spec.add(1, d + z * u * l + z * l**2)
9 spec.run_tuner(1)
```

$$\delta \geq \log(e^\zeta + e^{\zeta+\sigma})$$

$$\lambda \geq \log(e^\delta + e^{\zeta+\lambda} + e^{\zeta+2\lambda})$$

$$\lambda - 10,000\zeta \rightarrow \min_{\lambda, \sigma, \zeta}$$

# References

Boltzmann samplers compiler:

<https://github.com/maciej-bendkowski/generic-boltzmann-brain>