

Effective random generation of algebraic data types through Boltzmann samplers

Maciej Bendkowski

May 18th, 2023

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

where

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

Example: General purpose codecs

Let's start with implementing a *dictionary based* encoder for `String`, in the spirit of Lempel-Ziv-78.

Example: General purpose codecs

Encoding:

- ▶ initiate an alphabet $\Sigma = \{a, b, c, ..\}$
- ▶ scan the input $w_1 \dots w_n$ for its *longest prefix* $w_1 \dots w_p \in \Sigma$
- ▶ insert $w_1 \dots w_p w_{p+1}$ into Σ and output the position of $w_1 \dots w_p$ in Σ
- ▶ repeat until the whole input is encoded

```
1 | encode :: String -> [Int]
```

Example:

Σ	Input	Output
$\Sigma = \{a, b\}$	abaaa	0
$\Sigma = \{a, b, ab\}$	baaa	0 1
$\Sigma = \{a, b, ab, ba\}$	aaa	0 1 0
$\Sigma = \{a, b, ab, ba, aa\}$	aa	0 1 0 4

Example: General purpose codecs

Decoding:

- ▶ initiate an alphabet $\Sigma = \{a, b, c, ..\}$
- ▶ process pairs of successive code words (c_i, c_{i+1}) .
- ▶ for each pair, insert $c_i x$ into Σ where x is the initial letter of c_{i+1}
- ▶ output the c_i th symbol w_{c_i} in Σ
- ▶ repeat until the whole input is decoded

```
1 | decode :: [Int] -> String
```

Example:

Σ	Input	Output
$\Sigma = \{a, b\}$	0 1 0 4	a
$\Sigma = \{a, b, ab\}$	1 0 4	a b
$\Sigma = \{a, b, ab, ba\}$	0 4	a b a
$\Sigma = \{a, b, ab, ba, aa\}$	4	a b a aa

Example: General purpose codecs

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 = (...)
```

Example: General purpose codecs

```
1  instance Codec BinTree where
2      encode :: BinTree -> ByteString
3      encode = Binary.encode . LZ.encode . code
4
5      decode :: ByteString -> BinTree
6      decode = uncode . LZ.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
```

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              ]
```


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?*

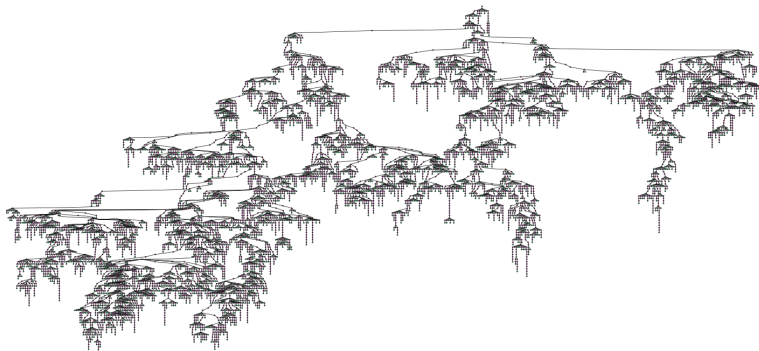
Example: General purpose codecs

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.

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

Example: Random λ -terms

Let's write a generator for λ -terms in the DeBruijn notation!



Example: Random λ -terms

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

Example: Random λ -terms

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

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
```

- ▶ 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`.

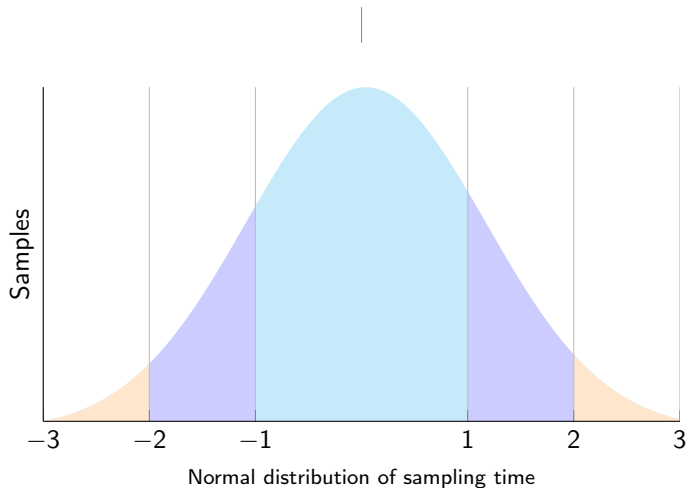
See <https://github.com/maciej-bendkowski/generic-boltzmann-brain>

Boltzmann samplers

```
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
```

Sampling efficiency

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



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

Beyond 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 }
```

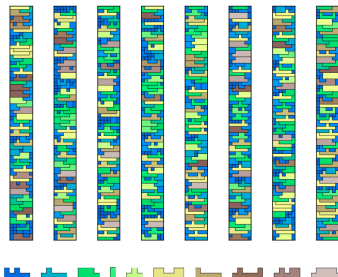


Figure: Five random $n \times 7$ tilings of areas in the interval $[500; 520]$ using in total 95 different tiles.

Univariate Boltzmann models

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

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

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 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

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

Univariate Boltzmann models

Boltzmann models admit elegant closure properties with respect to algebraic data type constructors, i.e. products and co-products (sums).

$$\{\omega\} \longrightarrow z^{|\omega|}$$

$$\mathcal{A} + \mathcal{B} \longrightarrow A(z) + B(z) = \sum_{n \geq 0} (a_n + b_n) z^n$$

$$\mathcal{A} \times \mathcal{B} \longrightarrow A(z) \times B(z) = \sum_{n \geq 0} \left(\sum_{k=0}^n a_k b_{n-k} \right) z^n$$

Example

The generating function $F(z)$ for Fibonacci numbers satisfies

$$F(z) = z + z^2 + (zF(z) + z^2F(z)) = \frac{z}{1 - z - z^2}$$

Univariate Boltzmann models

Given an algebraic data type \mathcal{S} and a target *mean size* n we can find the *branching probabilities* by solving

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

for the (unknown) x .

Problem

How can we solve this equation automatically, in the general case, of systems of possibly mutually recursive algebraic data types?

Convex optimisation: Random λ -terms

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

$$D(z) = z + zD(z)$$

$$L(z) = D(z) + zL(z) + zL(z)^2$$

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

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

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

Architecture

- ▶ `generic-boltzmann-brain` uses `paganini-hs` to formulate a *tuning problem* for `paganini` and then generates the sampler code at compile-time.
- ▶ `paganini`¹, a Python library, translates the domain-specific tuning problem into an optimisation problem using `cvxpy`, a modelling language for convex optimisation problems.
- ▶ `cvxpy`² uses `cvx` and a variety of optimisation solvers to find a solution to the original tuning problem.

¹<https://github.com/maciej-bendkowski/paganini>

²<https://www.cvxpy.org/>

Remarks

- ▶ **newtypes** are supported, and can be used if multiple samplers or size notions are simultaneously required.
- ▶ Polymorphic data types, such as

```
1 | data BinTree a
2 |   = Node (BinTree a) (BinTree a)
3 |   | Leaf a
```

are *not* supported. Since `a` is unknown at compile time, we cannot infer the *concrete* structure of `BinTree a` and formulate the corresponding tuning problem.

Conclusions

Thank you for your attention!

References:

- ▶ B. (2022) Automatic compile-time synthesis of entropy-optimal Boltzmann samplers
- ▶ B., Bodini, Dovgal (2018) Polynomial tuning of multiparametric combinatorial samplers
- ▶ Duchon, Flajolet, Louchard, Schaeffer (2004) Boltzmann Samplers for the Random Generation of Combinatorial Structures
- ▶ B. (2022) `generic-boltzmann-brain` @ github