

1. Coding

- *Huffman coding* (HC) efficient; optimally effective for bit-sequence-per-symbol
- *arithmetic coding* (AC)
 Shannon-optimal (fractional entropy); but computationally expensive
- asymmetric numeral systems (ANS) efficiency of Huffman, effectiveness of arithmetic coding
- applications of *streaming* (another story)

ANS introduced by Jarek Duda (2006–2013).

Now: Facebook (Zstandard), Apple (LZFSE), Google (Draco), Dropbox (DivANS)...

2. Intervals

Pairs of rationals

```
type Interval = (Rational, Rational)
```

with operations

```
unit = (0,1)

weight (l,r) x = l + (r - l) \times x

narrow i (p,q) = (weight i p, weight i q)

scale (l,r) x = (x - l) / (r - l)

widen i (p,q) = (scale i p, scale i q)
```

so that *narrow* and *unit* form a monoid, and inverse relationships:

```
weight i \ x \in i \iff x \in unit
weight i \ x = y \iff scale \ i \ y = x
narrow i \ j = k \iff widen \ i \ k = j
```

3. Models

```
Given
      counts::[(Symbol, Integer)]
get
      encodeSym:: Symbol → Interval
      decodeSym :: Rational → Symbol
such that
      decodeSym x = s \iff x \in encodeSym s
Eg alphabet {'a', 'b', 'c'} with counts 2, 3, 5 encoded as (0, \frac{1}{5}), (\frac{1}{5}, \frac{1}{2}), and (\frac{1}{2}, 1).
```

4. Arithmetic coding

```
encode_1 :: [Symbol] \rightarrow Rational
        encode_1 = pick \circ foldl \ estep_1 \ unit \ where
           estep_1 :: Interval \rightarrow Symbol \rightarrow Interval
           estep_1 is = narrow i (encodeSym s)
        decode_1 :: Rational \rightarrow [Symbol]
        decode_1 = unfoldr \ dstep_1 \ \mathbf{where}
           dstep_1 :: Rational \rightarrow Maybe (Symbol, Rational)
           dstep_1 \ x = let \ s = decodeSym \ x \ in \ Just \ (s, scale \ (encodeSym \ s) \ x)
where pick:: Interval \rightarrow Rational satisfies pick i \in I. Eg, with pick = fst:
        (0.1) \xrightarrow{\text{'a'}} (0.\frac{1}{5}) \xrightarrow{\text{'b'}} (\frac{1}{25}, \frac{1}{10}) \xrightarrow{\text{'c'}} (\frac{7}{100}, \frac{1}{10}) \sim \frac{7}{100}
```

5. Trading in bits

```
Let pick = fromBits \circ toBits :: Interval \rightarrow Rational, where toBits :: Interval \rightarrow [Bool] fromBits :: [Bool] \rightarrow Rational
```

Obvious thing: let *toBits i* pick shortest binary fraction in *i*, and *fromBits* evaluate this fraction. But this can't be streamed.

Instead, *toBits i* yields bit sequence *bs* such that bs + [True] is shortest:

```
toBits = unfoldr nextBit where

nextBit (l,r) \mid r \leq 1/2 = Just (False, (0, 1/2) \underline{widen} (l,r))

\mid 1/2 \leq l = Just (True, (1/2, 1) \underline{widen} (l,r))

\mid otherwise = Nothing

fromBits = foldr pack (1/2) where pack b x = ((if b then 1 else 0) + x) / 2
```

Now *pick* is a hylomorphism. Also, *toBits* yields a finite sequence.

6. Streaming encoding

Move *fromBits* from encoding to decoding:

```
encode_{Bits} :: [Symbol] \rightarrow [Bool]

encode_{Bits} = toBits \circ foldl \ estep_1 \ unit

decode_{Bits} :: [Bool] \rightarrow [Symbol]

decode_{Bits} = unfoldr \ dstep_1 \circ fromBits
```

Both of these can be *streamed*: alternately producing and consuming.

7. Towards ANS—fission and fusion

```
encode_1
          [[ definition; go back to pick = fst ]]
        fst ∘ foldl estep<sub>1</sub> unit
          [[ map fusion for foldl, backwards ]]
        fst o foldl narrow unit o map encodeSym
      = [[ narrow is associative ]]
        fst o foldr narrow unit o map encodeSym
      = [[ fusion for foldr ]]
        foldr weight 0 o map encodeSym
          [[ map fusion; let estep<sub>2</sub> s x = weight (encodeSym s) x ]]
        foldr estep<sub>2</sub> 0
so let encode_2 = foldr \ estep_2 \ 0.
```

8. Unfoldr-foldr theorem

Inverting a fold:

$$unfoldr\ g\ (foldr\ f\ e\ xs) = xs$$

$$\Leftarrow g(f \times z) = Just(x, z) \wedge ge = Nothing$$

Allowing *junk*:

$$(\exists ys. unfoldr \ g \ (foldr \ f \ e \ xs) = xs + ys) \iff g \ (f \ x \ z) = Just \ (x, z)$$

With *invariant*:

$$unfoldr g (foldr f e xs) = xs \qquad \iff ((g (f x z) = Just (x, z)) \iff p z) \land$$

 $((g e) = Nothing) \iff p e)$

where invariant *p* of *foldr f e* and *unfoldr g* is such that

$$p(f x z) \Leftarrow p z$$

 $p z' \Leftarrow p z \land g z = Just(x, z')$

9. Correctness of decoding

```
dstep_1 \ (estep_2 \ s \ z)

= [[ \ estep_2 \ ]]
dstep_1 \ (weight \ (encodeSym \ s) \ z)

= [[ \ dstep_1; \ let \ s' = decodeSym \ (weight \ (encodeSym \ s) \ z) \ ]]
Just \ (s', scale \ (encodeSym \ s') \ (weight \ (encodeSym \ s) \ z))

= [[ \ s' = s \ (see \ next \ slide) \ ]]
Just \ (s, scale \ (encodeSym \ s) \ (weight \ (encodeSym \ s) \ z))

= [[ \ scale \ i \circ weight \ i = id \ ]]
Just \ (s, z)
```

9. Correctness of decoding (continued)

```
Indeed, s' = s:
       decodeSym (weight (encodeSym s) z) = s
     weight (encodeSym s) z \in encodeSym s
     \iff [[ property of weight ]]
       z \in unit
and z \in unit is an invariant. Therefore
     take\ (length\ xs)\ (decode_1\ (encode_2\ xs)) = xs
for all finite xs.
```

10. From fractions to integers

AC encodes longer messages as more precise fractions. In contrast, ANS makes larger integers.

```
count :: Symbol \rightarrow Integer
cumul :: Symbol \rightarrow Integer
total :: Integer
find :: Integer \rightarrow Symbol

such that
```

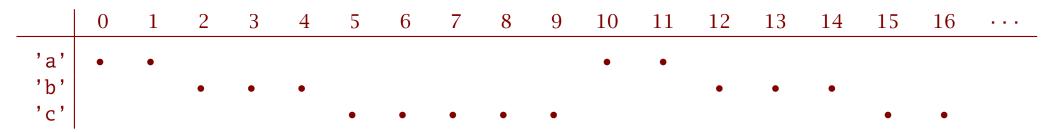
find $z = s \iff cumul \ s \leqslant z < cumul \ s + count \ s$ for $0 \leqslant z < total$.

11. Asymmetric encoding: the idea

- text encoded as integer z, with $\log_2 z$ bits of information
- next symbol s has probability p = count s / total, so requires $\log_2 (1/p)$ bits
- so map z, s to $z' \simeq z \times total / count <math>s$ —but do so invertibly
- with $z' = (z \underline{\text{div}} \ count \ s) \times total$, can undo the known multiplication: $z \underline{\text{div}} \ count \ s = z' \underline{\text{div}} \ total$
- what about s? headroom! with $z' = (z \underline{\text{div}} \ count \ s) \times total + cumul \ s$, $s = find (cumul \ s) = find (z' \underline{\text{mod}} \ total)$
- what about z? with $z' = (z \underline{\text{div}} \ count \ s) \times total + cumul \ s + z \underline{\text{mod}} \ count \ s$, $z \underline{\text{mod}} \ count \ s = z' \underline{\text{mod}} \ total cumul \ s$

12. ANS, pictorially

The start of the coding table for alphabet 'a', 'b', 'c' with counts 2, 3, 5:



- integers 0.. distributed across the alphabet, in proportion to counts
- encoding symbol s with current accumulator x yields the position of the xth blob in row s as the new accumulator x'

13. Essence of ANS encoding and decoding

```
encode_3 :: [Symbol] \rightarrow Integer
encode_3 = foldr \ estep_3 \ 0
estep_3 :: Symbol \rightarrow Integer \rightarrow Integer
estep_3 \ s \ z = let \ (q, r) = z \underline{divMod} \ count \ s \ in \ q \times total + cumul \ s + r
decode_3 :: Integer \rightarrow [Symbol]
decode_3 = unfoldr dstep_3
dstep_3 :: Integer \rightarrow Maybe (Symbol, Integer)
dstep_3 z = let (q, r) = z \underline{divMod} total
                   s = find r
              in Just (s, count s \times q + r - cumul s)
```

Correctness argument as before. But note that encoding is right-to-left.

14. Variation

Correctness does not depend on starting value: can pick any *l* instead of 0.

Also, $estep_3$ strictly increasing on z > 0, and $dstep_3$ strictly decreasing, so we know when to stop:

```
encode_4 :: [Symbol] \rightarrow Integer
encode_4 = foldr \ estep_3 \ l
decode_4 :: Integer \rightarrow [Symbol]
decode_4 = unfoldr \ dstep_4
dstep_4 :: Integer \rightarrow Maybe \ (Symbol, Integer)
dstep_4 \ z = \mathbf{if} \ z == l \ \mathbf{then} \ Nothing \ \mathbf{else} \ dstep_3 \ z
```

and we have

```
decode_4 (encode_4 xs) = xs
```

for all finite **xs**, but this time without junk.

15. Bounded precision

Fix base b and lower bound l. Represent accumulator z as pair (x, ys) such that:

- remainder ys is a list of digits in base b
- window x satisfies $l \le x < u$ for upper bound $u = l \times b$

under abstraction z = foldl inject x ys where

```
inject x y = x \times b + y and extract x = x \operatorname{divMod} b
```

Eg with b = 10 and l = 100, pair (123, [4, 5, 6]) represents 123456.

```
type Number = (Integer, [Integer])
```

Note "you can't miss it" properties:

inject
$$x y < u \iff x < l$$

 $l \le fst (extract x) \iff u \le x$

Want b, l powers of two, and u a single word. Also nice if $l \mod total = 0$.

16. Encoding

Maintain window in range.

```
econsume_5 :: [Symbol] \rightarrow Number
econsume_5 = foldr \ estep_5 \ (l, [])
estep<sub>5</sub> :: Symbol → Number → Number
estep_5 \ s \ (x, ys) = let \ (x', ys') = enorm_5 \ s \ (x, ys) \ in \ (estep_3 \ s \ x', ys')
enorm<sub>5</sub> :: Symbol → Number → Number
enorm_5 s(x, ys) = if estep_3 sx < u
                       then (x, ys)
                       else let (q, r) = extract x in enorm<sub>5</sub> s <math>(q, r : ys)
```

18

Pre-normalize *before* consuming symbol. Eg with b = 10, l = 100:

$$(340,[3]) \stackrel{\text{`a'}}{\leftarrow} (68,[3]) \stackrel{\text{norm}}{\leftarrow} (683,[]) \stackrel{\text{`b'}}{\leftarrow} (205,[]) \stackrel{\text{`c'}}{\leftarrow} (100,[])$$

17. Decoding

```
dproduce_5 :: Number \rightarrow [Symbol]
dproduce_5 = unfoldr \ dstep_5
dstep_5 :: Number \rightarrow Maybe \ (Symbol, Number)
dstep_5 \ (x, ys) =  let Just \ (s, x') = dstep_3 \ x
(x'', ys'') = dnorm_5 \ (x', ys)
in if x'' \geqslant l then Just \ (s, (x'', ys'')) else Nothing
dnorm_5 :: Number \rightarrow Number - dnorm_5 \ (enorm_5 \ s \ (x, ys)) = (x, ys) \ when \ l \leqslant x < u
dnorm_5 \ (x, y : ys) =  if x < l then dnorm_5 \ (inject \ x \ y, ys) else (x, y : ys)
dnorm_5 \ (x, []) = (x, [])
```

Decoding is symmetric to encoding: renormalize *after* emitting a symbol.

$$(340,[3]) \xrightarrow{\text{'a'}} (68,[3]) \xrightarrow{\text{norm}} (683,[]) \xrightarrow{\text{'b'}} (205,[]) \xrightarrow{\text{'c'}} (100,[])$$

Correctness again as before (no junk; invariant $l \le x < u$).

18. Trading in sequences

Flush out the last few digits in the *Number* when encoding complete:

```
eflush_5 :: Number \rightarrow [Integer]
eflush_5 (0, ys) = ys
eflush_5 (x, ys) = \mathbf{let} (x', y) = extract \ x \ \mathbf{in} \ eflush_5 (x', y : ys)
encode_5 :: [Symbol] \rightarrow [Integer]
encode_5 = eflush_5 \circ econsume_5
```

Conversely, populate initial *Number* from first few digits when decoding:

```
dstart_5 :: [Integer] \rightarrow Number

dstart_5 ys = dnorm_5 (0, ys)

decode_5 :: [Integer] \rightarrow [Symbol]

decode_5 = dproduce_5 \circ dstart_5
```

Note that $dstart_5$ ($eflush_5 x$) = $x \leftarrow l \leq x < u$.

19. Fast loops

```
encode :: [Symbol] \rightarrow [Integer] -- one tight loop
encode = h_1 l \circ reverse where
  h_1 \times (s:ss) = \text{let } x' = estep_3 \times x \text{ in if } x' < u \text{ then } h_1 \times x' \text{ ss else}
                     let (q, r) = extract \times in \ r : h_1 \ q \ (s : ss)
  h_1 x [] = h_2 x
  h_2 x = \text{if } x == 0 \text{ then } [] \text{ else let } (x', y) = extract x \text{ in } y : h_2 x'
decode :: [Integer] \rightarrow [Symbol] -- one tight loop
decode = h_0 \ 0 \circ reverse  where
  h_0 \ x \ (y : ys) \ | \ x < l = h_0 \ (inject \ x \ y) \ ys
  h_0 \times ys
              = h_1 \times ys
  h_1 \times ys = let Just(s, x') = dstep_3 \times in h_2 \cdot s \times ys
  h_2 \ s \ x \ (y : ys) \mid x < l = h_2 \ s \ (inject \ x \ y) \ ys
  h_2 s x ys = \mathbf{if} x \ge l \mathbf{then} s : h_1 x ys \mathbf{else}
```

20. But...

- AC encoding and decoding are both left-to-right
- so AC can be made *adaptive*: adjust model with each symbol
- fixed-precision AC is (I believe) equivalent to a quantizing adaptation

22

- ANS encoding and decoding go in *opposite* directions
- can't be so easily made adaptive
- then fixed-precision ANS is a different problem

Comments welcome! Paper in preparation...