

Deplump: A streaming lossless compressor

Nicholas Bartlett

Frank Wood

Department of Statistics, Columbia University, New York, USA

Abstract

1 Algorithm

Given an ordered symbol set Σ , probabilistic compression algorithms work by using a generative model to predict a sequence of symbols. The predictive distribution function is then used as the parameter in a range encoder to compress the stream. The details of a range encoder implementation are not included here, we only note that if the predictive distribution function is F and the next symbol in the stream is s then the parameters required by the range encoder are $F(s - 1)$ and $F(s)$. In the algorithm the function `RangeEncode()` implements the encoding and returns a bit sequence (possibly null). The use of a cumulative distribution function is well defined since the symbols are ordered. The notation $s - 1$ refers to the symbol prior to s in the symbol ordering. In order to decompress the stream the exact same predictive model will need to be built from the compressed stream. This requires that the model estimate prior to compressing s_n is a function of fixed parameters and the symbols $[s_0, s_1, \dots, s_{n-1}]$ because those are the only symbols available to the decompressor for decompressing s_n .

The algorithm operates primarily on a suffix tree. A suffix tree is a data structure for keeping track of unique suffixes of a set of strings. The tree structure arranges the suffixes hierarchically which makes it easy to search. In the case of a single stream the set of strings to consider is the set of contexts $\{[], [s_0], [s_0, s_1], [s_0, s_1, s_2], \dots\}$. Each node of the suffix tree corresponds to a context of the form $[s_m, \dots, s_{m+k}]$. In general we will use the notation \mathcal{N} to refer interchangeably to a node object and the context to which the node corresponds. The function `CreateNode(\mathcal{N}, \mathcal{M})` makes explicit the creation of a node \mathcal{N} associated with context \mathcal{N} and with parent \mathcal{M} . We use the function `PA(\mathcal{N})` to refer to the parent of node \mathcal{N} .

Each node object \mathcal{N} contains two integer counts for each $s \in \Sigma$, c_s, t_s . The notation c and t refer to the marginal counts in a node, namely $\sum_{s \in \Sigma} c_s$ and $\sum_{s \in \Sigma} t_s$. Each node also has a discount d associated with it. The discount associated with \mathcal{N} is a function of \mathcal{D} , the fixed discount parameters of the model, $|\mathcal{N}|$, and $|\text{PA}(\mathcal{N})|$. The function `GetDiscount` describes how the discount of a node \mathcal{N} is calculated. Suffix tree data structures use a reference sequence to maintain the unique suffixes in the tree. We use the notation \mathcal{RS} to refer to this reference sequence. Each node \mathcal{N} in a suffix tree contains two indices into \mathcal{RS} from which the context specific to node \mathcal{N} can be reconstructed. If the indices for \mathcal{N} are i and j , then the context associated with \mathcal{N} is $\mathcal{RS}[i : j]$.

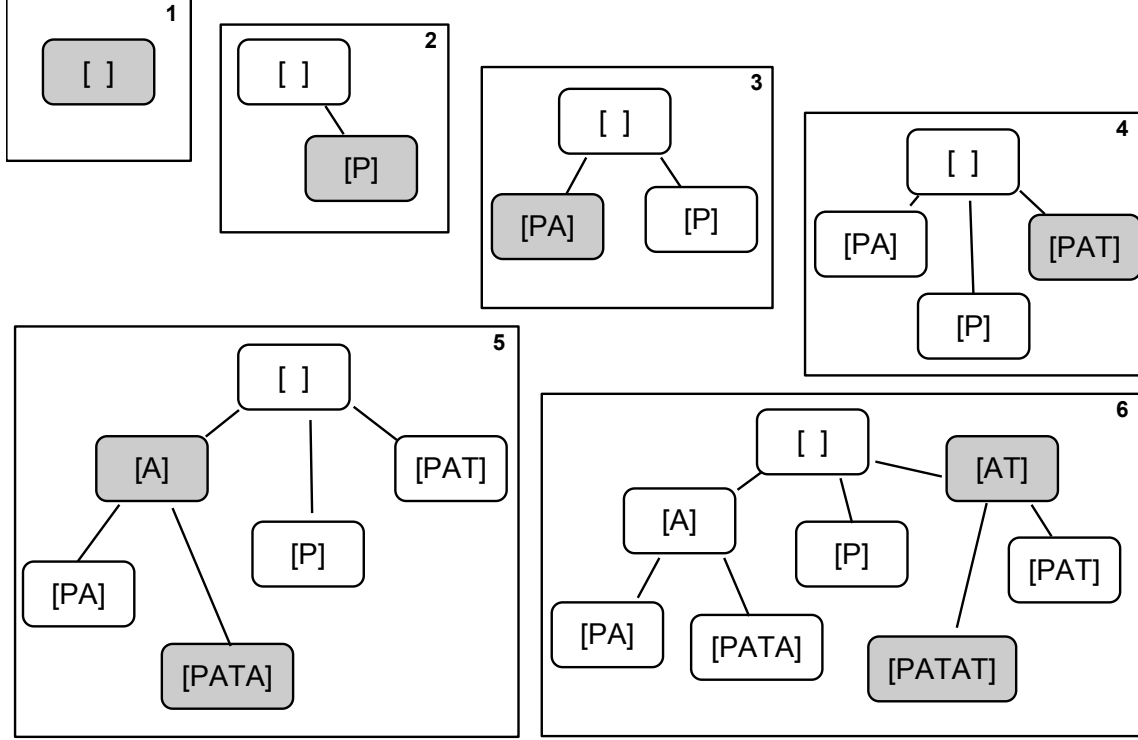


Figure 1: Construction of suffix tree for string “PATAT”. In each frame the new nodes are shaded in gray.

The reference sequence \mathcal{RS} is implemented as a linked list. Therefore, the indexing of a suffix tree node into \mathcal{RS} includes a pointer to a linked list node, an offset within that node and the length of the respective context. Each node of the linked list must also contain a list of suffix tree nodes which reference it. The reference sequence grows as the length of the compressed sequence grows and must be shortened as the algorithm progresses. The shortening of \mathcal{RS} is made explicit by the σ operator in `CDFNextSymbol`, which returns the input sequence shortened by removing the first element. When \mathcal{RS} is shortened, nodes in the suffix tree which reference removed sections are no longer usable and must be removed from the tree to prevent a memory leak. The cost of these operations can be amortized by shorting \mathcal{RS} in chunks by removing nodes of the linked list. Without a list of suffix tree nodes which reference each linked list node, deletion of the unusable nodes requires a search over the tree which is prohibitive for large trees. To minimize the impact of rendering nodes unusable by shortening the reference sequence, nodes should be updated to point to the most recent part of \mathcal{RS} as possible.

Since the model must be estimated incrementally, the suffix tree must also be incrementally constructed. Construction of the tree is handled by the function `GetNode` in Algorithm 1. An illustration of the incremental construction of a suffix tree can be seen in Figure 1 for the toy sequence [PATAT]. In frame 4 the function `GetNode` assigns [] to \mathcal{M} and then [PAT] to \mathcal{S} with $\mathcal{M} = \text{PA}(\mathcal{S})$. In Frame 5 `GetNode` assigns [PA] to \mathcal{M} , but then must assign [A] to \mathcal{P} with $\mathcal{P} = \text{PA}(\mathcal{M})$. \mathcal{S} is then created by `CreateNode([PATA], \mathcal{P})`. In

each frame the first step is to find \mathcal{M} , which can be achieved by descending an appropriate path of the suffix tree. All of the nodes on the path to \mathcal{M} and possibly \mathcal{M} itself can have the indices into \mathcal{RS} updated to point to a later section of the reference sequence. Finally, although not made explicit in the algorithm, it usually makes sense to limit the maximum length of a context in order to limit the depth of the suffix tree. In Section ?? we empirically investigate performance for three values of depth.

For each s in the input sequence the function `CDFNextSymbol` is used to obtain the predictive cumulative distribution function values of interest. After encoding the symbol first step to updating the model estimate is performed by the function `UpdateCountsAndDiscounts`. Starting at node \mathcal{N} and progressing up to the root of the tree, c_s is incremented if t_s was incremented in the node below. If c_s is incremented, a stochastic decision is made to increment t_s . The gradients for the discount parameters \mathcal{D} are updated by the function `UpdateDiscountParameterGradients`. Finally, if c is larger than k in any of the nodes, the counts c_s and potentially t_s are reduced by the function `ThinCounts`. The parameter k is a fixed parameter of the model.

First, c_s is incremented in node \mathcal{N} . A

to the root a stochastic decision is made to increment c_s

Things to make sure to include

Assumed ordering of the symbols in the symbol set. Describe parametrization of discounts
Mention learning rate η Describe `RangeEncode` function and return Describe σ operator
Describe size of tree L Describe length of \mathcal{RS} Describe the deletion of leaf nodes uniformly at random Describe `pa(\mathcal{N})` operator Describe count notation Describe `CreateNode` function
Mention arrays indexed from 0 Describe Draw multinomial Describe k in thin counts

Algorithm 1 Deplump

```
1: procedure DEPLUMP/PLUMP( $\mathcal{IS}, k$ )
2:    $\mathcal{RS} \leftarrow []$  ▷ reference sequence
3:   Initialize  $[]$  node of  $\mathcal{T}$  ▷ suffix tree
4:    $nc \leftarrow 1$  ▷ node count
5:    $\mathcal{D} \leftarrow \{d_0, d_1, d_2, \dots, d_{10}, \alpha\}$  ▷ discount parameters
6:    $\mathcal{G} \leftarrow \vec{0}$  ▷ discount parameter gradients,  $|\mathcal{G}| = |\mathcal{D}|$ 
7:    $\mathcal{OS} \leftarrow []$  ▷ output sequence
8:   if Plump then
9:      $h \leftarrow \text{InitializePointOnCDF}(\mathcal{IS})$ 
10:  end if
11:  for  $i = 0: |\mathcal{IS}|$  do
12:    if Plump then
13:       $[h, s, F(s-1), F(s)] \leftarrow \text{PointOnCDFNextSymbol}(h, \mathcal{T}, \mathcal{RS})$ 
14:       $b \leftarrow s$ 
15:    else
16:       $s = \mathcal{IS}[i]$ 
17:       $[F(s-1), F(s)] \leftarrow \text{CDFNextSymbol}(s_i, \mathcal{T}, \mathcal{RS})$ 
18:       $b \leftarrow \text{RangeEncode}(F(s-1), F(s))$ 
19:    end if
20:     $\mathcal{OS} \leftarrow [\mathcal{OS} \ b]$  ▷ append b to output sequence
21:     $\mathcal{D} \leftarrow \mathcal{D} + \mathcal{G}\eta / (F(s) - F(s-1))$  ▷ update discount parameters
22:     $\mathcal{G} \leftarrow \vec{0}$  ▷ reset gradients to zero
23:     $\mathcal{RS} \leftarrow [\mathcal{RS} \ s_i]$  ▷ append symbol to reference sequence
24:  end for
25:  return  $\mathcal{OS}$ 
26: end procedure
27: function CDFNEXTSYMBOL( $s, \mathcal{T}, \mathcal{RS}$ )
28:  while  $|\mathcal{RS}| \geq 100L$  do
29:    Delete nodes referencing  $\mathcal{RS}[0]$ 
30:     $\mathcal{RS} \leftarrow \sigma(\mathcal{RS})$ 
31:  end while
32:  while  $nc > (L - 2)$  do
33:    Delete leaf node uniformly at random
34:  end while
35:   $\mathcal{N} \leftarrow \text{GetNode}(\mathcal{RS}, \mathcal{T})$ 
36:   $\pi \leftarrow \text{PMF}(s, \mathcal{N}, \vec{0}, 1.0)$  ▷  $|\vec{0}| = |\Sigma|$ 
37:   $F(s-1) \leftarrow \sum_{i=1}^{s-1} \pi_i$ 
38:   $F(s) \leftarrow F(s-1) + \pi_s$ 
39:  UpdateCountsAndDiscountGradients( $\mathcal{N}, s, F(s) - F(s-1), \text{TRUE}$ )
40:  return  $[F(s-1), F(s)]$ 
41: end function
```

Algorithm 2 Deplump Continued

```
1: function POINTONCDFNEXTSYMBOL( $h, \mathcal{T}, \mathcal{RS}$ )
2:   while  $|\mathcal{RS}| \geq 100L$  do
3:     Delete nodes referencing  $\mathcal{RS}[0]$ 
4:      $\mathcal{RS} \leftarrow \sigma(\mathcal{RS})$ 
5:   end while
6:   while  $nc > (L - 2)$  do
7:     Delete leaf node uniformly at random
8:   end while
9:    $\mathcal{N} \leftarrow \text{GetNode}(\mathcal{RS}, \mathcal{T})$ 
10:   $\pi \leftarrow \text{PMF}(\mathcal{N}, \vec{0}, 1.0)$   $\triangleright |\vec{0}| = |\Sigma|$ 
11:   $F(s) \leftarrow 0$ 
12:   $s \leftarrow 0$ 
13:  while  $F(s) < h$  do
14:     $i \leftarrow i + 1$ 
15:     $F(s) \leftarrow F(s) + \pi_s$ 
16:  end while
17:   $F(s - 1) \leftarrow F(s) - \pi_s$ 
18:  Update  $h$  return  $[h, s, F(s - 1), F(s)]$ 
19: end function
20: function GETDISCOUNT( $\mathcal{N}$ )
21:   $d = 1.0$ 
22:  for  $i = (|\text{PA}(\mathcal{N})| + 1) : |\mathcal{N}|$  do
23:    if  $i \leq 10$  then
24:       $d \leftarrow d\mathcal{D}[i]$   $\triangleright$  multiply by discount parameter  $i$ 
25:    else
26:       $d \leftarrow d\mathcal{D}[10]^{\mathcal{D}[11]^i}$ 
27:    end if
28:  end for
29:  return  $d$ 
30: end function
31: function GETNODE( $\mathcal{S}, \mathcal{T}$ )
32:  Find the node  $\mathcal{M}$  in the suffix tree sharing the longest suffix with  $\mathcal{S}$ .
33:  if  $\mathcal{M}$  is a suffix of  $\mathcal{S}$  then
34:    if  $\mathcal{S} = \mathcal{M}$  then
35:      return  $\mathcal{M}$ 
36:    else
37:       $\mathcal{S} \leftarrow \text{CreateNode}(\mathcal{S}, \mathcal{M})$ 
38:      return  $\mathcal{S}$ 
39:    end if
40:  else
41:     $\mathcal{P} \leftarrow \text{FragmentNode}(\mathcal{M}, \mathcal{S})$ 
42:     $\text{PA}(\mathcal{M}) \leftarrow \mathcal{P}$ 
43:     $\mathcal{S} \leftarrow \text{CreateNode}(\mathcal{S}, \mathcal{P})$ 
44:    return  $\mathcal{S}$ 
45:  end if
46: end function
```

Algorithm 3 Deplump Continued

```

1: function UPDATECOUNTSANDDISCOUNTS( $\mathcal{N}$ ,  $s$ ,  $p$ , BackOff)
2:    $d \leftarrow \text{GetDiscount}(\mathcal{N})$ 
3:   if  $c > 0$  then
4:      $pp \leftarrow (p - \frac{c_i - t_i d}{c})(\frac{c}{td})$ 
5:      $tw \leftarrow c_i + d(t * pp - t_i)$ 
6:   else
7:      $pp \leftarrow p$ 
8:   end if
9:   if BackOff and  $c > 0$  then
10:     $c_s^{\mathcal{N}} \leftarrow c_s^{\mathcal{N}} + 1$ 
11:    BackOff  $\leftarrow 0$ 
12:    BackOff  $\leftarrow 1$  w.p.  $pp(\frac{t^{\mathcal{N}} d}{tw})$  ▷ w.p abbreviates “with probability”
13:    if BackOff then
14:       $t_s \leftarrow t_s + 1$ 
15:    end if
16:  else
17:    if BackOff then
18:       $c_s \leftarrow c_s + 1$ 
19:       $t_s \leftarrow t_s + 1$ 
20:    end if
21:  end if
22:  UpdateDiscountParameterGradients( $t_s$ ,  $t$ ,  $pp$ ,  $d$ )
23:  UpdateCountsAndDiscounts(PA( $\mathcal{N}$ ),  $s$ ,  $pp$ , BackOff)
24:  ThinCounts( $\mathcal{N}$ ,  $k$ )
25: end function
26: function THINCOUNTS( $\mathcal{N}$ ,  $k$ )
27:    $d \leftarrow \text{GetDiscount}(\mathcal{N})$ 
28:   while  $\sum_{s \in \Sigma} c_s > k$  do
29:      $s \leftarrow \text{DrawMultinomial}(\pi(\Sigma))$  s.t.  $\pi_l = \frac{c_l}{c}$ 
30:      $\phi \leftarrow \text{SamplePartition}(c_s, t_s, d)$ 
31:      $i \leftarrow \text{DrawMultinomial}([\frac{\phi_0}{c_s}, \frac{\phi_1}{c_s}, \dots, \frac{\phi_{t_s-1}}{c_s}])$ 
32:     if  $\phi_i == 1$  then
33:        $t_s \leftarrow t_s - 1$ 
34:     end if
35:      $c_s \leftarrow c_s - 1$ 
36:   end while
37: end function

```

Algorithm 4 Deplump Continued

```

1: function PMF( $\mathcal{N}, \pi, m$ )
2:    $d \leftarrow \text{GetDiscount}(\mathcal{N})$ 
3:   if  $c > 0$  then
4:     for  $i \in \Sigma$  do
5:        $\pi_i \leftarrow \pi_i + m(\frac{c_i - t_i d}{c})$ 
6:     end for
7:   end if
8:   if  $\text{PA}(\mathcal{N}) \neq \text{null}$  then
9:     return  $\text{PMF}(\text{PA}(\mathcal{N}), \pi, dm)$ 
10:  else
11:     $\pi \leftarrow (1 - dm)\pi + dm\mathcal{U}(\Sigma)$   $\triangleright \mathcal{U}(\Sigma)$  is the uniform distribution over  $\Sigma$ 
12:    return  $\pi$ 
13:  end if
14: end function
15: function FRAGMENTNODE( $\mathcal{M}, \mathcal{S}$ )
16:    $d^{\mathcal{M}} \leftarrow \text{GetDiscount}(\mathcal{M})$ 
17:    $\mathcal{P} \leftarrow$  maximum overlapping suffix of  $\mathcal{M}$  and  $\mathcal{S}$ 
18:    $\mathcal{P} \leftarrow \text{CreateNode}(\mathcal{P}, \text{PA}(\mathcal{M}))$ 
19:    $d^{\mathcal{P}} \leftarrow \text{GetDiscount}(\mathcal{P})$ 
20:   for  $s \in \Sigma$  do
21:      $\phi \leftarrow \text{SamplePartition}(c_s^{\mathcal{M}}, t_s^{\mathcal{M}}, d^{\mathcal{M}})$ 
22:      $c_s^{\mathcal{P}} \leftarrow 0$ 
23:      $t_s^{\mathcal{P}} \leftarrow t_s^{\mathcal{M}}$ 
24:      $t_s^{\mathcal{M}} \leftarrow 0$ 
25:     for  $i = 1 : |\phi|$  do
26:        $a \leftarrow \text{DrawCRP}(\phi[i], d^{\mathcal{M}}/d^{\mathcal{P}}, -d^{\mathcal{M}})$ 
27:        $t_s^{\mathcal{M}} \leftarrow t_s^{\mathcal{M}} + a$ 
28:        $c_s^{\mathcal{P}} \leftarrow c_s^{\mathcal{P}} + a$ 
29:     end for
30:   end for
31:   return  $\mathcal{P}$ 
32: end function

```

Algorithm 5 Deplump Continued

```
1: function SAMPLEPARTITION( $c, t, d$ )
2:    $M \leftarrow t \times c$  matrix of zeros
3:    $M(t, c) \leftarrow 1.0$ 
4:   for  $j = (c - 1) : 1$  do
5:     for  $i = 1 : (t - 1)$  do
6:        $M(i, j) \leftarrow M(i + 1, j + 1) + M(i, j + 1)(j - id)$ 
7:     end for
8:      $M(d, j) \leftarrow M(t, j + 1)$ 
9:   end for
10:   $\phi \leftarrow \vec{0}$   $\triangleright |\vec{0}| = t$ 
11:   $\phi[1] \leftarrow 1$ 
12:   $k \leftarrow 1$ 
13:  for  $j = 2 : c$  do
14:     $M(k, j) \leftarrow M(k, j)(j - 1 - kd)$ 
15:     $r \leftarrow 0$ 
16:     $r \leftarrow 1$  w.p.  $\frac{M(k+1, j)}{M(k+1, j) + M(k, j)}$ 
17:    if  $r = 1$  then
18:       $k \leftarrow k + 1$ 
19:       $\phi[k] \leftarrow 1$ 
20:    else
21:       $i \leftarrow \text{DrawMultinomial}([\frac{\phi[1]-d}{j-1-kd}, \frac{\phi[2]-d}{j-1-kd}, \dots, \frac{\phi[k]-d}{j-1-kd}])$ 
22:       $\phi[i] \leftarrow \phi[i] + 1$ 
23:    end if
24:  end for
25:  return  $\phi$ 
26: end function
27: function DRAWCRP( $n, d, c$ )
28:   $t \leftarrow 1$ 
29:  for  $i = 2 : n$  do
30:     $r \leftarrow 0$ 
31:     $r \leftarrow 1$  w.p.  $\frac{td+c}{i-1+c}$ 
32:     $t \leftarrow t + r$ 
33:  end for
34:  return  $t$ 
35: end function
36: function UPDATEDISCOUNTPARAMETERGRADIENTS( $p\_depth, depth, t_s, c, t, pp, d, m$ )
37:  if  $c > 0$  then
38:    if  $depth = 0$  then
39:       $\text{derivLogDd} \leftarrow \frac{1.0}{\mathcal{D}[\vec{0}]}$ 
40:       $\mathcal{G}[\vec{0}] \leftarrow \mathcal{G}[\vec{0}] + (d(t * pp - t_s)\text{derivLogDd}/c)m$ 
41:    else
42:       $z \leftarrow p\_depth + 1$ 
43:      while  $z \leq depth$  and  $z < 10$  do
44:         $\text{derivLogDd} \leftarrow \frac{1.0}{\mathcal{D}[z]}$ 
45:         $\mathcal{G}[z] \leftarrow \mathcal{G}[z] + (d(t * pp - \frac{t_s}{8})\text{derivLogDd}/c)m$ 
46:      end while
47:      if  $depth \geq 10$  then
48:         $a \leftarrow z - 10$ 
49:         $b \leftarrow depth - z + 1$ 
```