



Generalized Boundary Detection Using Compression-Based Analytics

Christina Ting, Richard Field, Tu-Thach Quach, and Travis Bauer

ABSTRACT

We present a new method for boundary detection within sequential data using compression-based analytics. Our approach is to approximate the information distance between two adjacent sliding windows within the sequence. Large values in the distance metric are indicative of boundary locations. A new algorithm is developed, referred to as sliding information distance (SLID), that provides a fast, accurate, and robust approximation to the normalized information distance. A modified smoothed z-score algorithm is used to locate peaks in the distance metric, indicating boundary locations. A variety of data sources are considered, including text and audio, to demonstrate the efficacy of our approach.

DATASETS

We used synthesized datasets with known ground truth in order to quantify algorithm performance:

1. Text randomly selected from English and Spanish translations of United Nations (UN) documents
2. Audio randomly selected from male and female speakers from the LibriSpeech ASR corpus

Each data set contains 100 synthesized boundaries with a section length of 512 tokens.

A. Sliding Information Distance (SLID)

- Approximate the information distance (ID) between adjacent subsequences x_k and y_k located at byte k in sequence z [Algorithm 1]

$$z_0 \dots z_{k-w-1} \underbrace{z_{k-w} \dots z_{k-1}}_{x_k} \underbrace{z_k \dots z_{k+w-1}}_{y_k} z_{k+w} \dots$$

- Form Lempel-Ziv (L - Z) dictionaries $D(x_k)$ and $D(y_k)$ [Algorithm 2]
- Compute S_k , the Jaccard distance between them at byte k

$$S_k(z; w) = 1 - \frac{|D(x_k) \cap D(y_k)|}{|D(x_k) \cup D(y_k)|}$$

- Slide the local window to the right and repeat the calculation to estimate this distance along z
 - Update to $D(x_{k+1})$ and $D(y_{k+1})$ by drop/append for efficiency [Algorithm 3]

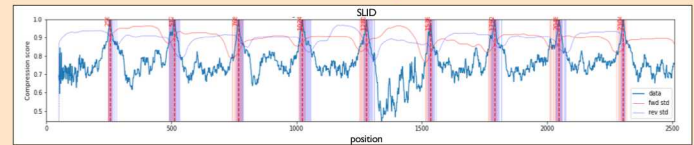
Algorithm 1 Sliding information distance.

```
1: function SLID(sequence  $z$ , window size  $w$ )
2:    $S \leftarrow [0]$ 
3:   for  $k = w, \dots, \text{do}$ 
4:      $x_k \leftarrow z_{k-w}, \dots, z_{k-1}$ 
5:      $y_k \leftarrow z_k, \dots, z_{k+w-1}$ 
6:     if  $k == w$  then
7:        $Lx, Dx \leftarrow \text{makeLZdict}(x_k)$ 
8:        $Ly, Dy \leftarrow \text{makeLZdict}(y_k)$ 
9:     else
10:       $Lx, Dx \leftarrow \text{updateLZdict}(x_k[-1], Lx)$ 
11:       $Ly, Dy \leftarrow \text{updateLZdict}(y_k[-1], Ly)$ 
12:    end if
13:     $S.append(1 - |Dx \cap Dy| / |Dx \cup Dy|)$ 
14:  end for
15:  return  $S$ 
16: end function
```

METHODS

B. Boundary Detection

- Any significant change in the SLID score indicates a possible boundary location
- We apply a smoothed z-score algorithm to identify peaks
 - The score is anomalous at k if it exceeds n standard deviations over the running mean of the previous m bytes
- We execute the z-score in forward and reverse to identify regions of local anomalous scores (shaded regions), then take the maximum



Algorithm 2 Initialize the LZ dictionary.

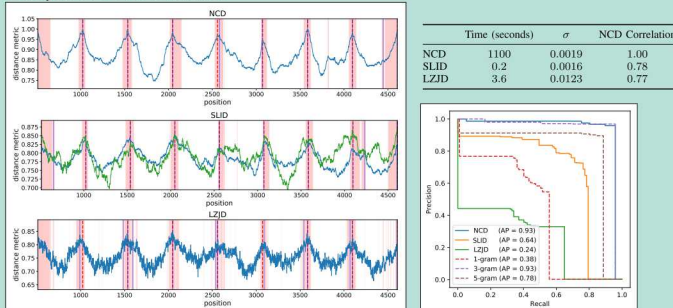
```
1: function MAKELZDICT(sequence  $b$ )
2:    $\ell \leftarrow []$ , start  $\leftarrow 0$ , end  $\leftarrow 0$ 
3:   while end <  $|b|$  do
4:     item  $\leftarrow b[\text{start} : \text{end}]$ 
5:     if item  $\notin \ell$  then
6:       start  $\leftarrow \text{end}$ 
7:     end if
8:      $\ell.append(\text{item})$ 
9:     end  $\leftarrow \text{end} + 1$ 
10:  end while
11:  return  $\ell$ , set( $\ell$ )
12: end function
```

Algorithm 3 Update the LZ dictionary.

```
1: function UPDTELZDICT(token  $t$ , list  $\ell$ )
2:   if  $\ell[-1] \neq t$  then
3:     item  $\leftarrow t$ 
4:   else
5:     item  $\leftarrow \ell[-1] + t$ 
6:   end if
7:    $\ell \leftarrow \ell[1:]$ 
8:    $\ell.append(\text{item})$ 
9:   return  $\ell$ , set( $\ell$ )
10: end function
```

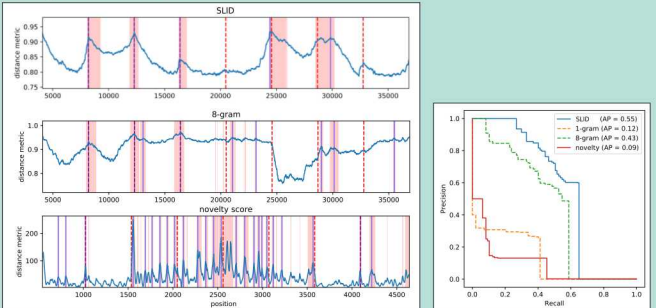
RESULTS

Analysis of UN Documents



- Compare SLID to Normalized Compression Distance (NCD), Lempel-Ziv Jaccard Distance (LZJD), n-Gram, and a Novelty Score (developed for audio signals)
- SLID is both accurate (see Precision-Recall curves) and efficient, while being robust to small changes in information distance
 - A noisy signal is more likely to produce false positives

Analysis of Audio Data



- If data is well-described by an n-gram, as is the case with text, a 3-gram performs best
 - Knowledge of underlying data needed to select optimal n
 - Compression is a good choice when there is no such knowledge
- SLID performs consistently well for both applications
 - Even outperforming the specialized novelty score for audio.

CONCLUSIONS

We have presented a fast, efficient, and robust compression-based method for detecting boundaries in arbitrary sequences of data, including data streams. We have demonstrated the versatility of our approach through several experiments on multiple data sources. As our method is computationally efficient, it is possible to apply different window sizes to obtain a multi-scale structural representation of the underlying data source for boundary detection. This could lead to further improved detection performance or provide the ability to identify coarse and fine-grained boundaries. We defer investigating this problem to our future work.