

Phrase Extraction Algorithm

Jakub Waszczuk

December 2018

1 Naive Algorithm

The naive phrase extraction algorithm (Alg. 1) considers all possible phrase pairs (\bar{e}, \bar{f}) given sentence pair (e, f) one by one. Each such phrase pair is added to the resulting set if it is consistent with the given alignment relation A . Recall that, graphically, each phrase pair (\bar{e}, \bar{f}) corresponds to a rectangle R in

Algorithm 1 Naive algorithm

```
 $F := \emptyset$ 
for each  $i, i': 1 \leq i \leq i' \leq n$  do
  for each  $j, j': 1 \leq j \leq j' \leq m$  do
    let  $\bar{f} = (f_i, \dots, f_{i'})$ 
    let  $\bar{e} = (e_j, \dots, e_{j'})$ 
    if  $(\bar{f}, \bar{e})$  consistent with  $A$  then
       $F := F \cup (\bar{f}, \bar{e})$ 
    end if
  end for
end for
```

the tabular representation of A , and that it is consistent if:

- R is non-empty (at least one marked cell) and
- for each column/row that intersects R , all its marked cells must be in R

The overall computation cost of the naive algorithm is $\Theta(m^2 \times n^2)$ times the cost of checking the consistency of a given rectangle (phrase pair) R w.r.t. A . The latter is not negligible, hence the total cost is unsatisfactory. There is a better way, fortunately.

2 Improved Algorithm

The improved algorithm (Alg. 2) considers all possible English phrases \bar{e} for a given sentence pair. For each such English phrase, it:

1. Identifies the minimal matching foreign phrase \bar{f}
2. Checks if (\bar{e}, \bar{f}) is consistent with A

Algorithm 2 Improved algorithm

```
 $F := \emptyset$ 
for each  $e_{\text{beg}}, e_{\text{end}} : 1 \leq e_{\text{beg}} \leq e_{\text{end}} \leq m$  do
  // Find the minimal matching foreign phrase
   $(f_{\text{beg}}, f_{\text{end}}) \leftarrow \text{MINIMALMATCHING}(e_{\text{beg}}, e_{\text{end}})$ 
  // Extract the phrase and its possible extensions
   $F \leftarrow F \cup \text{EXTRACT}(e_{\text{beg}}, e_{\text{end}}, f_{\text{beg}}, f_{\text{end}})$ 
end for
function MINIMALMATCHING( $e_{\text{beg}}, e_{\text{end}}$ )
  // Find the minimal matching foreign phrase
   $(f_{\text{beg}}, f_{\text{end}}) \leftarrow (n + 1, 0)$ 
  for each  $(e, f) \in A$  do
    if  $e_{\text{beg}} \leq e \leq e_{\text{end}}$  then
       $f_{\text{beg}} \leftarrow \min(f, f_{\text{beg}})$ 
       $f_{\text{end}} \leftarrow \max(f, f_{\text{end}})$ 
    end if
  end for
  return  $(f_{\text{beg}}, f_{\text{end}})$ 
end function
function EXTRACT( $e_{\text{beg}}, e_{\text{end}}, f_{\text{beg}}, f_{\text{end}}$ )
  // Check if at least one alignment point
  return  $\emptyset$  if  $f_{\text{end}} = 0$ 
  // Check if alignments points violate consistency
  for each  $(e, f) \in A$  do
    return  $\emptyset$  if  $f_{\text{beg}} \leq f \leq f_{\text{end}}$  and  $(e < e_{\text{beg}} \text{ or } e > e_{\text{end}})$ 
  end for
   $E \leftarrow \emptyset$ 
   $f_b = f_{\text{beg}}$ 
  repeat
     $f_e = f_{\text{end}}$ 
    repeat
      Add phrase  $(e_{\text{beg}} \dots e_{\text{end}}, f_b \dots f_e)$  to  $E$ 
       $f_e = f_e + 1$ 
    until  $f_e$  aligned
     $f_b = f_b - 1$ 
  until  $f_b$  aligned
  return  $E$ 
end function
```

3. If so, (\bar{e}, \bar{f}) is added to the resulting set, as well as all „extensions“ covering neighboring non-aligned words (both foreign and English)

This algorithm is faster than the naive one because, for a given \bar{e} , it doesn't consider all possible foreign phrases \bar{f} , only the relevant ones.

Proposition 1. *The improved algorithm 2 is not only faster than algorithm 1 but also correct – both calculate the same set of phrase pairs.*

Let $x_{(i,j)}$ be the phrase spanning (i, j) in sentence \mathbf{x} . To prove the above proposition we need to show that, if $(\bar{e}, f_{(f_{\text{beg}}, f_{\text{end}})})$ is consistent with A , then:

1. The minimal matching phrase $f_{(f_{\text{beg}}, f_{\text{end}})}$ is really minimal – i.e., for any foreign span (i, j) such that either $i > f_{\text{beg}}$ or $j < f_{\text{end}}$, $(\bar{e}, f_{(i,j)})$ is not consistent with A .
2. Extending $f_{(f_{\text{beg}}, f_{\text{end}})}$ with neighboring non-aligned words (both foreign and English) leads to a consistent phrase pair.
3. Extending $f_{(f_{\text{beg}}, f_{\text{end}})}$ with neighboring aligned words leads, again, to a inconsistency.

We now proceed to show the first property. Let $i > f_{\text{beg}}$ (the proof is analogous in case $j < f_{\text{end}}$). From the definition of MINIMALMATCH, we can see that $(e, f_{\text{beg}}) \in A$ for some $e : e_{\text{beg}} \leq e \leq e_{\text{end}}$. Since $e_{\text{beg}} \leq e \leq e_{\text{end}}$ and $(e, f_{\text{beg}}) \in A$, consistency requires that $i \leq f_{\text{beg}} \leq j$. However, that contradicts the initial assumption that $i > f_{\text{beg}}$.