# A Survey of Phrase Projectivity in Antigone

Jonathan Sterling

April 2013

## 1   Dependency Trees and Their Projectivity

A dependency tree encodes the head-dependent relation for a string of words, where arcs are drawn from heads to their dependents. We consider a phrase *projective* when these arcs do not cross each other, and *discontinuous* to the extent that any of the arcs intersect. Figure 1 illustrates the various kinds of projectivity violations that may occur.



μεστῇ πολλῶν ἀγαθῶν

πολλῶν μεστὸν ἀγαθῶν

(a) "Full of plentiful supplies" (Xenophon, *Anabasis* 3.5.1) is fully projective.

(b) "Full of many good things" (Plato, *Laws* 906a) has one projectivity violation.

στὰς δ' ὑπὲρ μελάθρων φονώσαισιν ἀμφιχανὼν κύκλῳ λόγχαις ἑπτάπυλον στόμα ἔβα

(c) "And he stood over the rooftops, gaped in a circle with murderous spears around the seven-gated mouth, and left" (Sophocles, *Antigone* 117–120) has five projectivity violations (note that multiple arcs may intersect at a point).
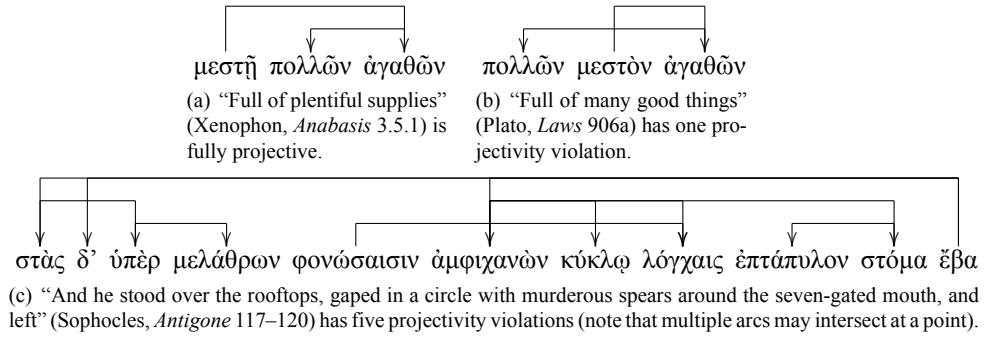
Figure 1: A dependency path wrapping around itself is a projectivity violation, as in (b); interlacing adjacent phrases also violate projectivity, as in (c). Examples (a–b) drawn from Devine & Stephens.

In this paper, we use a concrete metric of projectivity $\omega$, given by the following ratio:

$$\omega = \frac{\text{number of violations}}{\text{number of arcs}}$$

Section 2 deals with the development of an algorithm to compute this quantity for a particular dependency tree.

## 2   Algorithm & Data Representation

Dependency trees are a recursive data structure with a head node, which may have any number of arcs drawn to further trees (this is called a *rose tree*). We represent them as a Haskell data-type as follows:

$$\textbf{data } \mathsf{Tree}\ \alpha = \alpha \curvearrowright [\mathsf{Tree}\ \alpha]$$

This can be read as "For all types $\alpha$, a Tree of $\alpha$ is constructed from a *label* of type $\alpha$ and a *subforest* of Trees of $\alpha$," where brackets are a notation for lists.

Given a tree, we can extract its root label or its subforest by pattern matching on its structure as follows:

```
getLabel :: Tree α → α
getLabel (l ⌢ _) = l
getForest :: Tree α → [Tree α]
getForest (_ ⌢ ts) = ts
```

## 2.1  From Edges to Trees

A sentence from the Perseus treebank is in the form of a list of words that are indexed by their linear position, and cross-referenced by the linear position of their dominating head. We shall consider each index to be a *vertex*, and each pair of vertices to be an Edge, which we shall write as follows:

$$\textbf{data } \mathsf{Edge}\ \alpha = \alpha \leftrightarrow \alpha\ \textbf{deriving } \mathsf{Eq}$$

An Edge $\alpha$ is given by two vertices of type $\alpha$; the **deriving** Eq statement generates the code that is necessary to determine whether or not two Edges are equal using the ($\equiv$) operator. In order to perform our analysis, we should wish to transform the raw list of edges into a tree structure. The basic procedure is as follows:

First, we try to find the root vertex of the tree. This will be a vertex that is given as the head of one of the words, but does not itself appear in the sentence:

```
rootVertex :: Eq α ⇒ [Edge α] → Maybe α
rootVertex es = find (∉ deps) heads where
    heads = ⟦ (λ(x ↔ y) → x) es ⟧
    deps  = ⟦ (λ(x ↔ y) → y) es ⟧
```

If the data that we are working with are not well-formed, there is a chance that we will not find a root vertex; that is why the type is given as Maybe.

Then, given a root vertex, we look to find all the edges that it touches, and try to build the subtrees that are connected with those edges.

```
onEdge :: Eq α ⇒ α → Edge α → 𝔹
onEdge i (x ↔ y) = x ≡ i ∨ y ≡ i

oppositeVertex :: Eq α ⇒ α → Edge α → α
oppositeVertex i (x ↔ y)
    | x ≡ i      = y
    | otherwise = x
```

This is done recursively until the list of edges is exhausted and we have a complete tree structure:

```
treeFromEdges :: Ord α ⇒ [Edge α] → Maybe (Tree α)
treeFromEdges es = ⟦ (buildWithRoot es) (rootVertex es) ⟧ where
   buildWithRoot es root = root ↷ sortedChildren where
      roots          = ⟦ (oppositeVertex root) localVertices ⟧
      children       = ⟦ (buildWithRoot foreignVertices) roots ⟧
      localVertices  = filter (onEdge root) es
      foreignVertices = filter (¬ ∘ onEdge root) es
      sortedChildren = sortBy (compare 'on' getLabel) children
```

## 2.2 Counting Violations: Computing $\omega$

Violations are given as an integer tally:

```
type Violations = ℤ
```

The basic procedure for counting projectivity violations is as follows: flatten down the tree into a list of edges cross-referenced by their vertical position in the tree; then traverse the list and see how many times these edges intersect each other.

```
type Level = ℤ
```

The vertical position of a node in a tree is represented as its Level, counting backwards from the total depth of the tree. That is, the deepest node in the tree is at level $0$, and the highest node in the tree is at level $n$, where $n$ is the tree's depth.

```
levels :: Tree α → [[α]]
levels t = fmap (fmap getLabel) $
   takeWhile (¬ ∘ null) $
   iterate (≫=getForest) [t]


depth :: Tree α → ℤ
depth = length ∘ levels
```

We can now annotate each node in a tree with what level it is at:

```
annotateLevels :: Tree α → Tree (Level, α)
annotateLevels tree = aux (depth tree) tree where
   aux l (x ↷ ts) = (l, x) ↷ ⟦ (aux (l − 1)) ts ⟧
```

Then, we fold up the tree into a list of edges and levels:

```
allEdges :: Ord α ⇒ Tree α → [(Level, Edge α)]
allEdges tree = aux (annotateLevels tree) where
   aux ((_, x) ↷ ts) = ts ≫= go where
      go t@((l, y) ↷ _) = (l, edgeWithRange [x, y]) : aux t
```

3

```
edgeWithRange :: (Ord α) ⇒ [α] → Edge α
edgeWithRange xs = minimum xs ↔ maximum xs
```

A handy way to think of edges annotated by levels is as a representation of the arc itself, where the vertices of the edge are the endpoints, and the level is the height of the arc.

   If one end of an arc is between the ends of another, then there is a single intersection. If one arc is higher than another and the latter is in between the endpoints of the former, there is no violation; but if they are at the same level, or if the latter is higher than the former, there is a double intersection. Otherwise, there is no intersection.

```
checkEdges :: Ord α ⇒ (Level, Edge α) → (Level, Edge α) → Violations
checkEdges (l, xy@(x ↔ y)) (l', uv@(u ↔ v))
   | y ∈_E uv ∧ u ∈_E xy            = 1
   | x ∈_E uv ∧ v ∈_E xy            = 1
   | x ∈_E uv ∧ y ∈_E uv ∧ l ⩾ l' = 2
   | u ∈_E xy ∧ v ∈_E xy ∧ l ⩽ l' = 2
   | otherwise                      = 0
```

We determine whether a vertex is in the bounds of an edge using $\cdot \in_E \cdot$.

```
· ∈_E · :: Ord α ⇒ α → Edge α → 𝔹
z ∈_E x ↔ y = z > minimum [x, y]
              ∧ z < maximum [x, y]
```

We can now use what we've built to count the intersections that occur in a collection of edges. This is done by adding up the result of checkEdges of the combination of each edge with the subset of edges which are at or below its level:

```
edgeViolations :: Ord α ⇒ [(Level, Edge α)] → Violations
edgeViolations xs = sum ⟦ violationsWith xs ⟧ where
   rangesBelow (l, _) = filter (λ(l', _) → l' ⩽ l) xs
   violationsWith x   = sum ⟦ (checkEdges x) (rangesBelow x) ⟧
```

Finally, $\omega$ is computed for a tree as follows:

```
ω :: Ord α ⇒ Tree α → ℚ
ω tree = violationsCount / totalArcsCount  where
   edges           = allEdges tree
   violationsCount = fromIntegral (edgeViolations edges)
   totalArcsCount  = length edges
```

# 3   Parsing the Perseus Treebank

The Persues treebank is a collection of XML files, which have data in the following (simplified) scheme:

```
<sentence id="2900759">
  <word id="1" form="χρὴ" lemma="χρή" head="0" />
  <word id="2" form="δὲ" lemma="δέ" head="1" />
  ...
</sentence>

<sentence id="2900760">
  <word id="1" form="μεγάλοι" lemma="μέγας" head="3" />
  <word id="2" form="δὲ" lemma="δέ" head="12" />
  ...
</sentence>
```

We can express the general shape of such a document as follows:

$$
\begin{aligned}
&\textbf{newtype } \mathsf{XML} &&= \mathsf{XML}\ [\mathsf{Content}] \\
&\textbf{newtype } \mathsf{Word} &&= \mathsf{Word\ Element} \\
&\textbf{newtype } \mathsf{Sentence} &&= \mathsf{Sentence\ Element}
\end{aligned}
$$

To convert XML into trees, we must first extract the sentences from the file, and then we convert those into trees.

$$
\begin{aligned}
&\mathsf{sentencesFromXML} :: \mathsf{XML} \rightarrow [\mathsf{Sentence}] \\
&\mathsf{sentencesFromXML}\ (\mathsf{XML}\ xml) = \textbf{do} \\
&\quad elems \leftarrow \mathsf{onlyElems}\ xml \\
&\quad [\![\ \mathsf{Sentence}\ (\mathsf{findElements}\ (\mathsf{simpleName}\ \texttt{"sentence"})\ elems)\ ]\!]
\end{aligned}
$$

To build a tree from a sentence, first we get all of the words from that sentence and convert them into edges.

$$
\begin{aligned}
&\mathsf{wordsFromSentence} :: \mathsf{Sentence} \rightarrow [\mathsf{Word}] \\
&\mathsf{wordsFromSentence}\ (\mathsf{Sentence}\ s) = [\![\ \mathsf{Word}\ (\mathsf{findChildren}\ (\mathsf{simpleName}\ \texttt{"word"})\ s)\ ]\!]
\end{aligned}
$$

Edges are the content of the `head` attribute paired with that of the `id` attribute.

$$
\begin{aligned}
&\mathsf{edgeFromWord} :: \mathsf{Word} \rightarrow \mathsf{Maybe}\ (\mathsf{Edge}\ \mathbb{Z}) \\
&\mathsf{edgeFromWord}\ (\mathsf{Word}\ w) = [\![\ (\mathsf{readAttr}\ \texttt{"head"}\ w)\ \leftrightarrow\ (\mathsf{readAttr}\ \texttt{"id"}\ w)\ ]\!]
\end{aligned}
$$

Thence, we can build a tree from a sentence.

$$
\begin{aligned}
&\mathsf{treeFromSentence} :: \mathsf{Sentence} \rightarrow \mathsf{Maybe}\ (\mathsf{Tree}\ \mathbb{Z}) \\
&\mathsf{treeFromSentence} = \mathsf{treeFromEdges} \circ \mathsf{edgesFromSentence}\ \textbf{where} \\
&\quad \mathsf{edgesFromSentence} :: \mathsf{Sentence} \rightarrow [\mathsf{Edge}\ \mathbb{Z}] \\
&\quad \mathsf{edgesFromSentence}\ s = \mathsf{catMaybes}\ [\![\ \mathsf{edgeFromWord}\ (\mathsf{wordsFromSentence}\ s)\ ]\!]\ \textbf{where}
\end{aligned}
$$

By putting the pieces together, we also derive a function to read all the trees from an XML document:

> treesFromXML :: XML → [Tree ℤ]
> treesFromXML $xml$ = catMaybes ⟦ treeFromSentence (sentencesFromXML $xml$) ⟧

Finally, we must read the file as a string, parse it as XML, and then convert that XML into a series of trees.

> treesFromFile :: FilePath → IO [Tree ℤ]
> treesFromFile $path$ = ⟦ (treesFromXML ∘ XML ∘ parseXML) (readFile $path$) ⟧

# 4    Analysis of Data

We compute the average $\omega$ of the trees contained in a file as follows:

> analyzeFile :: FilePath → IO ℚ
> analyzeFile $path$ = **do**
>     $trees$ ← treesFromFile $path$
>     return (average ⟦ $\omega$ $trees$ ⟧)

# Appendix: Auxiliary Functions

> simpleName :: String → QName
> simpleName $s$ = QName $s$ Nothing Nothing

> readAttr :: Read $\alpha$ ⇒ String → Element → Maybe $\alpha$
> readAttr $n$ = fmap read ∘ findAttr (simpleName $n$)

> average :: Fractional $n$ ⇒ $[n]$ → $n$
> average $xs$ = $\dfrac{\text{sum } xs}{\text{length } xs}$