## Extracting Phrases with Chunking

This is a simple chunker based on an averaged perceptron. Here, accuracy is not key since we have a method that's tolerable to noise. And furthermore, a chunk in and of itself is completely useless (at least for my purposes). The chunks here are used to generate

1. Possible phrases of interest and

2. Factoring sentences into noun phrase verb phrase noun phrase.

The accuracy of 2) is not paramount, it simply needs to be good enough to feed a later hypothesis generation phase. A human will then be tasked to filter the outputs. If the signal is above the noise then that's all that matters.

```
1:  let formChunks chunks =
2:    let _, final_chunk, built_chunks =
3:     chunks |> Array.fold (fun ((lasttag, curphrase, phrases) as state) (w, pos) ->
4:           match splitstrWith "-" pos, curphrase with
5:             | [|_; chunktag|], _  ->
6:               if lasttag = chunktag then chunktag, (w,chunktag)::curphrase, phrases
7:               else chunktag, [w,chunktag], curphrase::phrases
8:             | _ , [] -> state
9:             | _ , c -> "", [], c::phrases) ("", [],[])
10:    final_chunk :: built_chunks |> List.rev |> List.toArray
```

The above method works to simply join a series of matching tags (for tagged words) together. E.g. "NP"; "NP"; "VP" => "NP" "NP"; "VP"... etc.

```
1:  let formChunksAny chunks =
2:    let _, final_chunk, built_chunks =
3:     chunks |> Array.fold (fun ((lasttag, curphrase, phrases) as state) (w, chunktag) ->
4:               if lasttag = chunktag then chunktag, (w,chunktag)::curphrase, phrases
5:               else chunktag, [w,chunktag], curphrase::phrases) ("", [],[])
6:
7:    final_chunk :: built_chunks |> List.rev |> List.toArray
```

This method is identical to the last but further removes the restriction that the word must be tagged.

### Joining Chunks.

```
1:   let joinchunksToStringGen f phraseList =
2:       phraseList |> Array.mapFilter (List.rev >> f >> joinTokens) (String.length >> (<) 1)
```

Joins chunks to string, dropping tags and any word whose len is $< 1$.
*Sample Output:*

```
1:   [|[|"Henry"; "was"; "an early king"; "of"; "France"|];
2:     [|"the boy"; "will throw"; "the ball"|]...
3:
```

joinchunksToStringTags keeps tags.

```
1:   ///phrase list as input
2:   let joinChunksToString phraseList = joinchunksToStringGen (List.map fst) phraseList
3:
4:   ///phrase list as input
5:   let joinChunksToStringTags phraseList = joinchunksToStringGen id phraseList
```

*Sample Output*:

```
1:   [|[|"(Henry, NP)"; "(was, VP)"; "(an, NP) (early, NP) (king, NP)";
2:     "(of, PP)"; "(France, NP)"|];
3:     [|"(the, NP) (boy, NP)"; "(will, VP) (throw, VP)"; "(the, NP) (ball, NP)"|];...
```

```
1:   let joinChunksToSummaryGen innermap phraseList =
2:       phraseList |> Array.filterMap (function [_, "NP"] | [_, "VP"] -> true | [] | [_,_] -> fal
3:                                     (List.rev >> innermap)
4:
5:   let joinChunksToSummary joinChar phrases = joinChunksToSummaryGen (List.map fst >> joinToken
```

### Example usage:

```
1:   let chunks = sentWords |> Array.Parallel.mapi (fun i s ->
2:         let t = predictlabelChunker tagtableChunk avec_chunk pos.[i] s
3:         Array.zip s t)
4:
5:   let gs = chunks |> Array.map getFocusPhrases
6:   let qq = gs.[3] |> Array.filter (List.isEmpty >> not)
7:   let qs = gs |> Array.filter (fun x -> x |> Array.filter (List.isEmpty >> not) |> Array.lengt
8:   qs.[2]
9:   let groupedChunks = chunks |> Array.map formChunks
10:
11:  let sections =  groupedChunks |> Array.map (joinChunksToSummary " ")
12:
13:  let allchunks = groupedChunks
14:                    |> Array.concat
15:                    |> Array.filterMap
```

```
16:                         (List.isEmpty >> not)
17:                         (fun phraseList ->
18:                                 phraseList
19:                                 |> List.rev
20:                                 |> List.map fst
21:                                 |> joinTokens,
22:                                     snd phraseList.Head) // |> Seq.g
23:
24:
25:   let chunk_tagmap = Map allchunks
26:
27:   let chunk_tag = allchunks |> Seq.groupBy snd |> Seq.mapGroupByWith (Set.ofSeq >> Set.toArray
28:   let chunk_map = Map chunk_tag
29:
30:   chunk_tag |> Array.map (keepLeft (Array.sortByDescending String.length))
31:
32:   chunk_map.["ADJP"] |> Array.sortByDescending (String.length)
33:
34:   sections |> Array.map (breakParagraphs splitwSpace 3 9 "\n" >> snd >> trim)
35:
36:   let paras = sections |> set |> Set.toArray
```

A more lenient grouping of chunks. Allows:

- verbs after nouns

- nouns after verbs.

- adverbs after verbs

- adjp after verb

```
1:   let groupChunkFlexible allowadverb_verb chunk =
2:       let _, last_chunk, built_chunks =
3:           chunk |> Array.fold (fun ((lasttag, curphrase, phrases) as state) (w, pos) ->
4:                   match splitstrWith "-" pos, curphrase with
5:                   | [|_; chunktag|], _   ->
6:                       if lasttag = chunktag || (chunktag = "VP" && lasttag = "NP") || (chunktag = '
7:                           || (allowadverb_verb && (chunktag = "ADJP" || chunktag = "ADVP") && (lastt
8:                       then chunktag, (w,chunktag)::curphrase, phrases
9:                       else chunktag, [w,chunktag], curphrase::phrases
10:                  | _ , [] -> state
11:                  | _ , c -> "", [], c::phrases) ("", [],[])
12:       last_chunk::built_chunks |> List.rev |> List.toArray
13:
14:
15:   let chunks0 = chunks |> Array.map (groupChunkFlexible true >> joinChunksToSummary ", ")
16:
```

```
17:   Array.zip (chunks0 |> Array.map (Seq.filter ((<>) ' ') >> Seq.length)) (sections |> Array.ma
18:
19:   let chunks9 = chunks |> Array.map (groupChunkFlexible true
20:                                >> joinChunksToSummaryGen (fun phrases ->
21:                                        phrases |> List.map fst
22:                                                |>  joinTokens,
23:                                    Seq.mode (phrases |> List.map snd)))
24:
25:   let paras0 = chunks9 |> Array.concat |> set |> Set.toArray
```

4