# Averaged perceptron tagger

Frankie Robertson
`frrobert@student.jyu.fi`

University of Jyväskylä

`https://www.github.com/frankier/perceptron-tagger-slides/`

18th of July, 2016

- Averaged structured perceptron with beam search
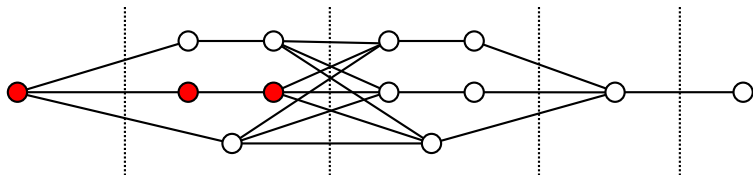- A DSL for features
- Results for Kazakh

- ▶ **Structured prediction** is breaking down prediction of a large structure into subproblems and using results from earlier stages in later stages (like dynamic programming). For POS tagging the simplest approach is break the problem down to per-token tagging and then tag left to right.

- ▶ **The perceptron** stores a sparse linear vector of (feature, weight) pairs. Observations are scored by taking the dot product of their feature vector with the perceptron. Training is done trying a prediction, and updating incorrect weights by reenforcing or penalising them depending on whether they correspond with the correct observation.

- ▶ With the **structured perceptron** the score from each update is accumulated to get the score of the (partial) output.

- ▶ The perceptron only converges for linearly separable inputs, otherwise its weights oscilate. To get the best weights the perceptron is trained for several iterations on reshuffled observations and the **averaged** weights across the whole training period are saved as the final weights.

- ▶ With **beam search**, the n-best candidates are considered and updated at each stage as opposed to a pure-greedy strategy which would only keep one intermediate result.

# Eurovision     ән     конкурсы 2010     .



Eurovision<np><al><nom>

    ән<n><nom> +e<cop><aor><p3><pl>
    ән<n><nom> +e<cop><aor><p3><sg>
    ән<n><nom>

        конкурс<n><px3sp><nom>
        конкурс<n><px3sp><nom> +e<cop><aor><p3><pl>
        конкурс<n><px3sp><nom> +e<cop><aor><p3><sg>

          2010<num><subst><nom>

           .<sent>

# A DSL for features 1/3

```
spectie: How's it going?

frankier: Been coding like a maniac

frankier: I've ended up creating a sort of lisp in XML
based on a stack VM

spectie: O___O
```

# A DSL for features 3/3

- ▶ Greenspun's tenth rule of programming: Any sufficiently complicated C or Fortran program contains an ad hoc, informally-specified, bug-ridden, slow implementation of half of Common Lisp.
- ▶ Neither easy nor desirable to try and define every possible useful feature up front.
- ▶ Seems likely that any feature language will be subject to feature creep (pun intended) so an expression language is a reasonable starting point.
- ▶ Bytecode interpreter design is faster and more modular than a tree walking interpreter and is usually the same or less code. Also it's easy to serialise.
- ▶ Hopefully the flexibility empowers language authors more than it confuses them. (In the worst case people can just copy-paste.)

# Example of DSL 1/2

```
<def-global as=''major_tag_0''>
  <subscript idx=''0''>
    <ex-tags>
      <ex-wordoid><wrdaddr /></ex-wordoid>
    </ex-tags>
  </subscript>
</def-global>

<def-global as=''is_dmorph''>
  <streq val=''+''>
    <slice end=''1''>
      <ex-lemma>
        <ex-wordoid><wrdaddr /></ex-wordoid>
      </ex-lemma>
    </slice>
  </streq>
</def-global>
```

# Example of DSL 2/2

```
<feat>
  <pred><var name=''is_headword'' /></pred>
  <out><var name=''lemma_0'' /></out>
  <out><var name=''major_tag_0'' /></out>
</feat>

<feat>
  <pred><var name=''is_dmorph'' /></pred>
  <out><var name=''headword_major_tag_0'' /></out>
  <out><var name=''major_tag_0'' /></out>
</feat>
```

# Results for Kazakh

Results!

# Future work

- Extend globals to macros/templates.
- Constructor for scanning/selecting a previous or subsequent wordoid or surface form or wordoid based on a predicate, eg to get the previous verb.
- Allow easy generation of multiple prefix/postfix features with special construct

# References

★ *Syntactic Processing using the Generalized Perceptron and Beam Search*. Zhang & Clark 2011. Includes general formulation and multiple worked examples. Most of the implementation is based on this paper.

▶ *Discriminative Training Methods for Hidden Markov Models*. Michael Collins 2002. Original Collins perceptron paper.

▶ *A Good Part-of-Speech Tagger in about 200 Lines of Python*, Matthew Honnibal (blog post). Plus: Easily understandable reference implementation. Minus: Formulation varies from rest of the literature.