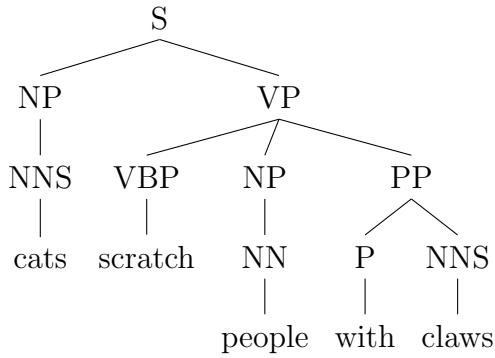


Integer Linear Programming for Semantic Role Labeling

Jacob Eisenstein

October 23, 2014

Here's our example:



Let's say that the correct annotation is

- $[cats]_{a0}$
- $[scratch]_v$
- $[people]_{a1}$
- $[with\ claws]_{am-mnr}$
- $[with]_{\emptyset}$
- $[claws]_{\emptyset}$
- $[scratch\ people\ with\ claws]_{\emptyset}$

(The tag AM-MNR means **argument modifier – manner**)

For each potential argument i , and for each tag t , define

$$Y_{i,t} = \begin{cases} 1, & \text{argument } i \text{ takes tag } t \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

Now suppose we define a feature function, where

$$\mathbf{f}(y, \mathbf{w}, i, t) = \delta(Y_{i,t} = y), \quad (2)$$

and set the weight

$$\theta(y, \mathbf{w}, i, t) = \log P(\mathbf{w}_i | y_i = t). \quad (3)$$

Then we have,

$$\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}, \mathbf{w}) = \sum_{i,t} \boldsymbol{\theta}^\top \mathbf{f}(y_{i,t}, \mathbf{w}, i, t) \quad (4)$$

$$= \sum_i \log P(\mathbf{w}_i | y_i = t) \quad (5)$$

$$\approx \log P(\mathbf{w} | \mathbf{y}, \tau), \quad (6)$$

where τ is the parse tree.

More reasonably, we can define a rich local feature function $\mathbf{f}(y_{i,t}, \mathbf{w}, i, t)$, using features like the ones given by Gildea and Manning. In either case, J is a **linear** function of the integer (binary) variables \mathbf{y} . Such optimization problems are called **Integer Linear Programs** (ILP). They fit into our usual form for linear prediction:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{w}, \tau)} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}, \mathbf{w}, \tau), \quad (7)$$

where $\mathcal{Y}(\mathbf{w}, \tau)$ is the set of allowable labelings.

Because \mathbf{y} is decoupled, we can optimize this easily: just solve separately for each y_i , obtaining the best label t for each constituent i . But this could lead to some bad outcomes.

- Multiple arguments of the same type:

[cats]_{a0} scratch [people]_{a0} with claws

- Overlapping arguments:

[with [claws]_{a1}]_{a2}

- Illegal arguments for the predicate:

[cats]_{a0} scratch [people]_{a1} [with claws]_{a2}¹

- Reference arguments without referents:

– ok: *[The deregulation]_{A1} of railroads [that]_{R-A1} began [in 1980]_{AM-TMP}*

– not ok: *[The deregulation]_{A1} of railroads [that]_{R-A0} began [in 1980]_{AM-TMP}*

Rather than incorporating global factors into the objective, we'll treat them as **constraints**, and solve a constrained optimization problem,

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{w}, \tau)} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}, \mathbf{w}, \tau) \quad (8)$$

$$s.t. \mathbf{y} \in \mathcal{C}(\mathbf{w}, \tau) \quad (9)$$

Where the constraint set $\mathcal{C}(\mathbf{w}, \tau)$ requires things like:

- All arguments get at most one label, $\forall i \sum_t y_{i,t} = 1$. Note we use equality, because you can always have the \emptyset label.

¹I'm not sure if a2 is really illegal here, but the point is that scratch is not a ditransitive verb, so it can't take a second direct object.

- No duplicate argument classes, $\forall t \neq \emptyset, \sum_i y_{i,t} \leq 1$
- Overlapping arguments get at most one non-null label:

$$\forall \langle i, j \rangle : i \rightsquigarrow_{\tau} j, y_{i,\emptyset} + y_{j,\emptyset} \geq 1 \quad (10)$$

- Some arguments are forbidden, e.g. $\sum_i y_{i,A2} = 0$
- Reference arguments must have referents

$$\begin{aligned} \forall i y_{i,R-A0} &\leq \sum_j y_{j,A0} \\ \forall i y_{i,R-A1} &\leq \sum_j y_{j,A1}, \dots \end{aligned}$$

- Continuation argument must follow the referent

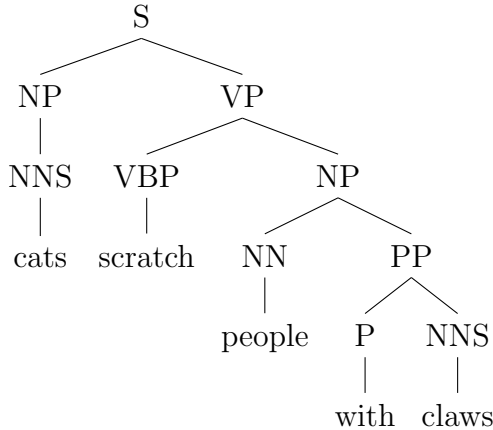
$$\begin{aligned} \forall i y_{i,C-A0} &\leq \sum_{j < i} y_{j,A0} \\ \forall i y_{i,C-A1} &\leq \sum_{j < i} y_{j,A1}, \dots \end{aligned}$$

All of the constraints are linear, meaning we can write them as $\mathbf{A}\mathbf{y} \leq \mathbf{b}$.

Unfortunately, constrained integer linear programming (ILP) is NP-Hard. It's easy to see this by reduction from TSP, boolean satisfiability, or minimum vertex cover. However, there are well-optimized algorithms for ILP, such as cutting plane and branch and bound, implemented in solvers like GLPK (free) or CPLEX (expensive). You just precompute the scores $\psi_{i,t} = \boldsymbol{\theta}^T \mathbf{f}(y_{i,t}, \mathbf{w}, i, t)$ for all $\langle i, t \rangle$, write out the constraints, and send it to the solver.

Punyakanok et al (2004) "It only takes about 10 minutes to solve the inference problem for all 4305 sentences on a Pentium-III 800 Mhz machine in our experiments."

Adding more parses Now, suppose we ran another parser and got another parse:



Now we have another possible argument: [people with claws]. We can just include this argument with the ones that we got from the other parse, as long as we add the appropriate constraints to prevent overlapping.

This allows us to combine many possible parse trees, and even to select arguments that are constituents from different trees. Punyakanok, Roth, and Yi (2007) showed that this flexibility yields a big boost in accuracy.