

Dependency grammars

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

- ▶ Based on syntactic relations between a **head** word and a **dependent** word
- ▶ “Shallower” in hierarchical structure

Add WeChat powcoder

Head and dependents

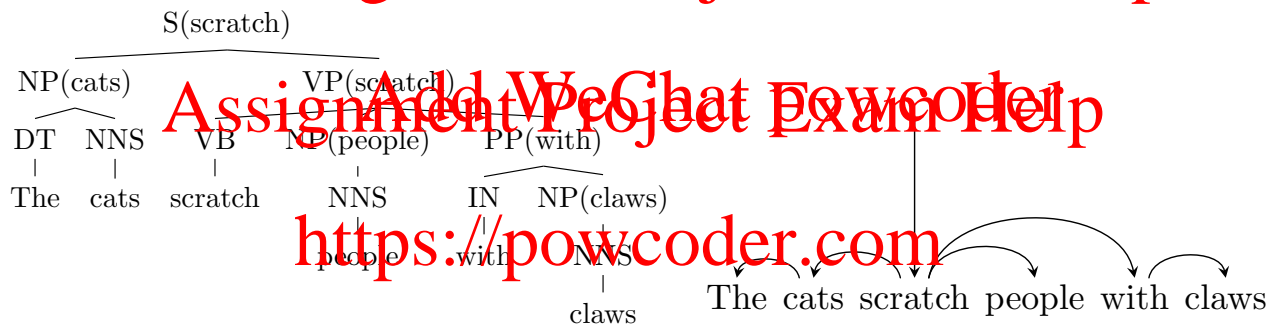
<https://powcoder.com>

- ▶ The head sets the syntactic category of the construction: for example, nouns are the heads of noun phrases, and verbs are the heads of verb phrases.
- ▶ The modifier may be optional while the head is mandatory: for example, in the sentence *Cats scratch people with claws*, the subtrees *cats scratch* and *cats scratch people* are grammatical sentences, but *with claws* is not.
- ▶ The head determines the morphological form of the modifier: for example, in languages that require gender agreement, the gender of the noun determines the gender of the adjectives and determiners.
- ▶ Edges should first connect content words, and then connect function words.

Relationship between phrase structures and dependency structures

<https://powcoder.com>

Assignment Project Exam Help



(a) lexicalized consistency parse (b) unlabeled dependency tree

Figure 11.1: Dependency grammar is closely linked to lexicalized context free grammars: each lexical head has a dependency path to every other word in the constituent. (This example is based on the lexicalization rules from § 10.5.2, which make the preposition the head of a prepositional phrase. In the more contemporary Universal Dependencies annotations, the head of *with claws* would be *claws*, so there would be an edge *scratch* → *claws*.)

Labeled dependencies

<https://powcoder.com>

Assignment Project Exam Help

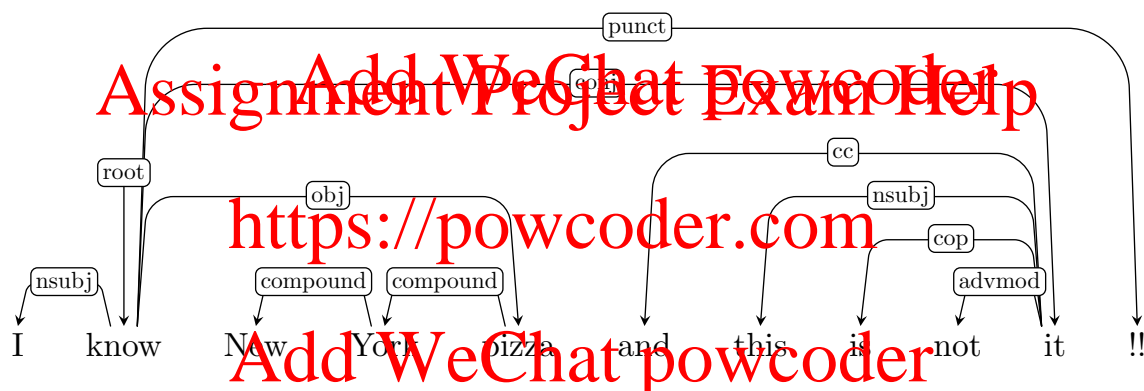


Figure 11.3: A labeled dependency parse from the English UD Treebank (reviews-361348-0006)

Projectivity

<https://powcoder.com>

- **Projectivity:** An edge from i to j is projective iff all k between i and j are descendants of i . A dependency parse is projective iff all its edges are projective.
- Informally, a dependency parse is projective if there are no crossing edges if all dependencies are drawn on one side of the sentence.

<https://powcoder.com>

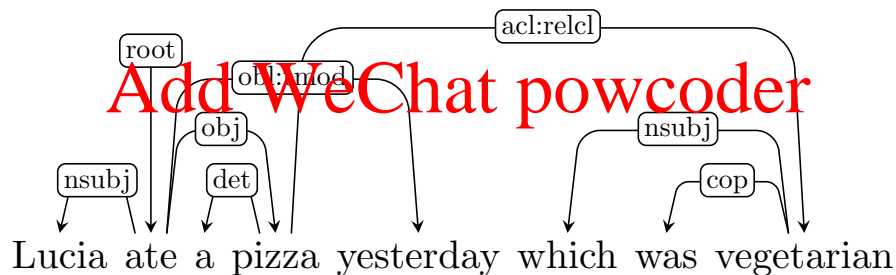


Figure 11.5: An example of a non-projective dependency parse. The “crossing edge” arises from the relative clause *which was vegetarian* and the oblique temporal modifier *yesterday*.

Main dependency parsing approaches

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

- ▶ Graph-based approaches
- ▶ Transition-based approaches

<https://powcoder.com>

Add WeChat powcoder

Graph-based approach

- ▶ Let $\mathbf{y} = \{(i \xrightarrow{r} j)\}$ represent a dependency graph in which r is a relation from headword $i \in \{1, 2, \dots, M, ROOT\}$ to modifier $j \in \{1, 2, \dots, M\}$. The special node $ROOT$ indicates the root of the graph, and M is the length of the input $|\mathbf{w}|$.
- ▶ Given a scoring function $\Psi(\mathbf{y}, \mathbf{w}; \theta)$,

$$\hat{\mathbf{y}} = \underset{\mathbf{y} \in \mathcal{Y}(\mathbf{w})}{\operatorname{argmax}} \Psi(\mathbf{y}, \mathbf{w}; \theta)$$

Add WeChat powcoder

where $\mathcal{Y}(\mathbf{w})$ is the set of valid dependency parses on the input \mathbf{w} .

- ▶ The set of possible labels $|\mathcal{Y}(\mathbf{w})|$ is exponential in the length of the input.
- ▶ Algorithm that search over this space of possible graphs are known as **graph-based dependency parsers**.

Factorization

Dependency parsers that operate under this assumption are called arc-factored.

- First-order factorization:

$$\Psi(\mathbf{y}, \mathbf{w}; \theta) = \sum_{i \xrightarrow{r} j} \psi(i \xrightarrow{r} j, \mathbf{w}; \theta)$$

- Second-order factorization:

$$\begin{aligned} \Psi(\mathbf{y}, \mathbf{w}; \theta) = & \sum_{i \xrightarrow{r} j} \psi(i \xrightarrow{r} j, \mathbf{w}; \theta) \\ & + \sum_{k \xrightarrow{r'} i \in \mathbf{y}} \psi_{grandparent}(i \xrightarrow{r} j, k, r', \mathbf{w}; \theta) \\ & + \sum_{i \xrightarrow{r'} s \in \mathbf{y}, s \neq j} \psi_{sibling}(i \xrightarrow{r} j, s, r', \mathbf{w}; \theta) \end{aligned}$$

Computing scores for dependency arcs

<https://powcoder.com>

Assignment Project Exam Help

- Assignment Project Exam Help
- ▶ Linear $\psi(i \rightarrow j, \mathbf{w}; \theta) = \theta \cdot \mathbf{f}(i \rightarrow j, \mathbf{w})$
 - ▶ Neural $\psi(i \xrightarrow{r} j, \mathbf{w}; \theta) = \text{Feedforward}([\mathbf{u}_w; \mathbf{u}_{w_j}]; \theta)$
 - ▶ Generative: $\psi(i \xrightarrow{r} j, \mathbf{w}; \theta) = \log P(w_j, r | w_i)$

<https://powcoder.com>
Add WeChat powcoder

Linear feature-based arc scores

<https://powcoder.com>

Assignment Project Exam Help

- ▶ The length and direction of the arc;
- ▶ The words w_i and w_j linked by the dependency relation;
- ▶ The prefixes, suffixes, and parts-of-speech of these words;
- ▶ The neighbors of the dependency arc, w_{i-1} , w_{i+1} , w_{j-1} , w_{j+1} ;
- ▶ The prefixes, suffixes, and part-of-speech of these neighbor words.

Learning a linear model with Perceptron

<https://powcoder.com>

- For a model with feature-based arc scores and perceptron loss, we obtain the usual structured perceptron update
 - Finding the tree with the highest score

<https://powcoder.com>

$$\hat{y} = \operatorname{argmax}_{y' \in \mathcal{Y}(\mathbf{w})} \theta \cdot \mathbf{f}(\mathbf{w}, y')$$

- Update the weights

$$\theta = \theta + \mathbf{f}(\mathbf{w}, y) - \mathbf{f}(\mathbf{w}, \hat{y})$$

Learning a linear model with CRF

<https://powcoder.com>

- ▶ A CRF for arc-factored dependency parsing is built on the probability model:

$$p(\mathbf{y}|\mathbf{w}) = \frac{\exp \sum_{i \rightarrow j \in \mathbf{y}} \psi(i \xrightarrow{r} j, \mathbf{w}; \theta)}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp \sum_{i \rightarrow j \in \mathbf{y}'} \psi(i \xrightarrow{r} j, \mathbf{w}; \theta)}$$

Where the posterior represents the fraction of the exponent of the score of one possible dependency graph out of the sum of the exponent of the scores of all possible graphs

- ▶ Questions: How do we compute the score of one dependency graph? How do we compute the sum of scores for all possible graphs?
- ▶ Such a model is trained to minimize the negative log conditional-likelihood.

Neural arc scores

- <https://powcoder.com>
- Assignment Project Exam Help
- Add WeChat powcoder
- ▶ Given vector representations of \mathbf{x}_i for each word w_i in the input, a set of arc scores can be computed from a feedforward neural network.

$$\psi(i \xrightarrow{r} j) = \text{FeedForward}([\mathbf{x}_i, \mathbf{x}_j], \Theta_r)$$

where unique weights Θ_r are available for each arc type.

- <https://powcoder.com>
- Add WeChat powcoder
- ▶ Specifically, the score of each arc for each relation can be computed as:

$$\mathbf{z} = g(\Theta_r[\mathbf{x}_i; \mathbf{x}_j] + b_r^{(z)})$$

$$\psi(i \xrightarrow{r} j) = \beta_r \mathbf{z} + b_r^{(y)}$$

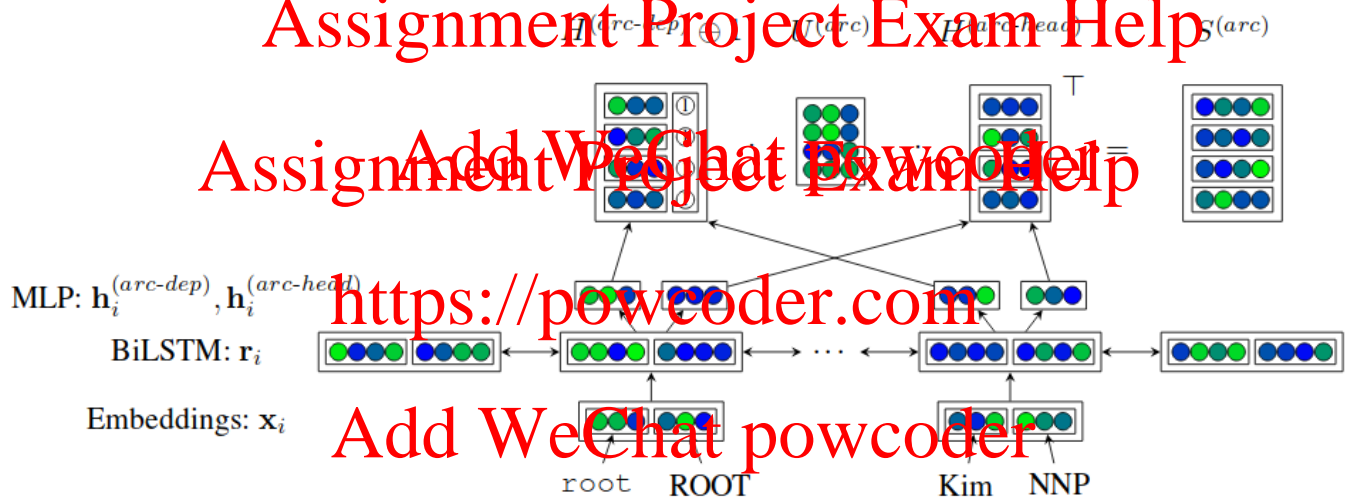
where Θ_r is a matrix, β_r is a vector, each b_r is a scalar, the function g is an element-wise *tanh* activation function.

SOA in neural dependency parsing: Biaffine models

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help



Tutorial: <http://www.cse.chalmers.se/~richajo/nlp2019/I7/Biaffine%20dependency%20parsing.html>

Bi-affine models for dependency parsing

- ▶ The input is a sequence of word embeddings concatenated with their POS embeddings

$$\mathbf{x}_i = \mathbf{v}_i^{(word)} \oplus \mathbf{v}_i^{(pos)}$$

- ▶ The input is fed into a BiLSTM to get a sequence of hidden states

$$\mathbf{h}_i = \text{BiLSTM}(\mathbf{x}_1, \dots, \mathbf{x}_n)_i$$

- ▶ Each hidden state \mathbf{h}_i is projected into four distinct vectors

$$\mathbf{h}_i^{(arc-dep)} = \text{MLP}^{(arc-dep)}(\mathbf{h}_i)$$

$$\mathbf{h}_i^{(arc-head)} = \text{MLP}^{(arc-head)}(\mathbf{h}_i)$$

$$\mathbf{h}_i^{(rel-dep)} = \text{MLP}^{(rel-dep)}(\mathbf{h}_i)$$

$$\mathbf{h}_i^{(rel-head)} = \text{MLP}^{(rel-head)}(\mathbf{h}_i)$$

Predicting the arcs with hidden vectors

<https://powcoder.com>

Assignment Project Exam Help

- ▶ Given a dependent, pair it up with each potential head (all other tokens) in the sentence, and compute a score:

$$s_i^{(arc)} = H^{(arc-head)} W^{(arc)} h_i^{(arc-dep)} + H^{(arc-head)} \mathbf{b}^T (arc)$$

<https://powcoder.com>

- ▶ The head is the pair with the highest score

Add WeChat powcoder

$$y_i^{(arc)} = \operatorname{argmax}_j s_{ij}^{(arc)}$$

Predicting relations with hidden states

- ▶ Given a head y_i' for word i , perform another biaffine transformation to predict the relation labels
- ▶ First use the relation vectors to compute a score for each possible label:

$$s_i^{(rel)} = \mathbf{h}_{y_i'}^{(rel-head)\top} \mathbf{U}^{(rel)} \mathbf{h}_i^{(rel-dep)} + \mathbf{W}^{(rel)} \left(\mathbf{h}_i^{(rel-dep)} \oplus \mathbf{h}_{y_i'}^{(rel-head)} \right) + \mathbf{b}^{(rel)}$$

What's the shape of \mathbf{U} ?

- ▶ Find the relation label with the highest score:

$$y_i'^{(rel)} = \operatorname{argmax}_j s_{ij}^{(rel)}$$

Breaking potential cycles in factored graph-based dependency parsing

<https://powcoder.com>

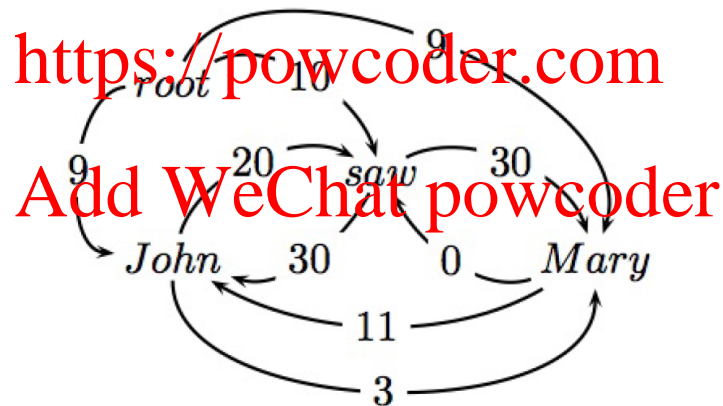
Assignment Project Exam Help

- ▶ Since in factored graph-based dependency parsing the parent for each dependent is predicted independently by finding the edge with the highest score, it is possible that the resulting dependency structure may have cycles. For instance, the model might predict i is the head of j , and j is the head of i .
- ▶ When this happens, we need to break the cycle to ensure that the resulting dependency tree is well-formed.
- ▶ One algorithm that can do this is the Chu-Liu/Edmonds algorithm.

The Chu-Liu-Edmonds algorithm

<https://powcoder.com>

- Assuming a model that assigns a score to each possible edge in a dependency tree
- $x = \text{root John saw Mary}$

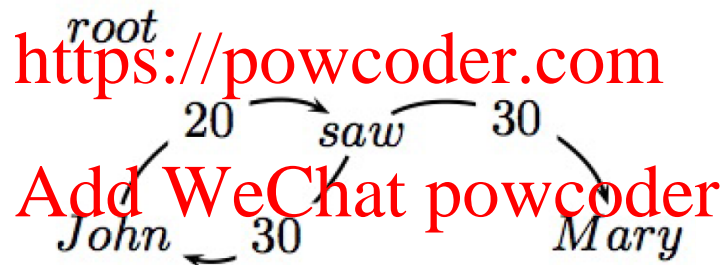


The Chu-Liu-Edmonds algorithm

<https://powcoder.com>

Assignment Project Exam Help

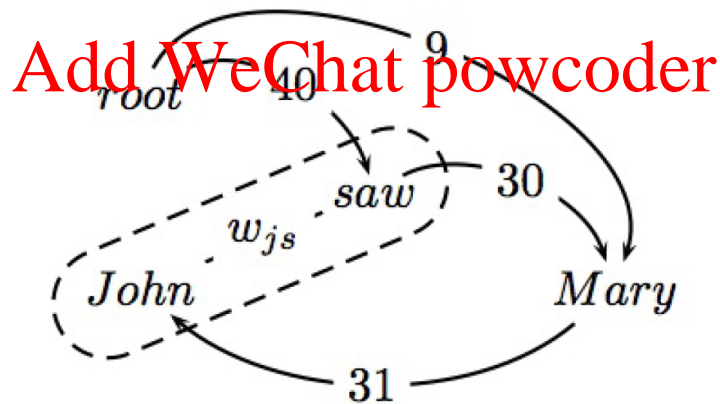
- ▶ start by removing all incoming arcs to *root*, and finding the highest scoring incoming arc for each node



The Chu-Liu-Edmonds algorithm

<https://powcoder.com>

- ▶ If not a tree, identify cycles and contract
- ▶ Recalculate arc weights into and out of cycle
- ▶ New incoming arc weights equal to the weight of the best spanning tree that includes the head of the incoming arc, and all nodes in cycle
- ▶ $\text{root} \rightarrow \text{saw} \rightarrow \text{John}: 40$
- ▶ $\text{root} \rightarrow \text{John} \rightarrow \text{saw}: 40$

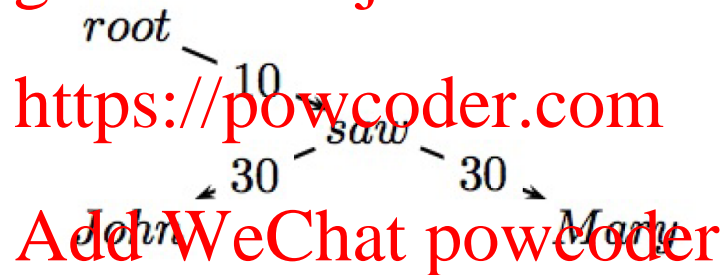


The Chu-Liu-Edmonds algorithm

<https://powcoder.com>

Assignment Project Exam Help

- This is the final dependency tree



Add WeChat powcoder