

KU Leuven - NLP 2021-2022
23 November 2021

Exercise Session 3

Syntactic Parsing and Semantic Role Labelling with GNN

For questions contact:
[Toledo Forum for Questions](#) or
nlp@ls.kuleuven.be

In this session we discuss syntactic constituency parsing and dependency parsing as well as Semantic Role Labelling (SRL). Background can be found in chapters 12, 13, and 14 of “Speech and Language Processing” by Dan Jurafsky and James H. Martin (2019, 3rd edition) of which a draft can be found online at: <https://web.stanford.edu/~jurafsky/slp3/>

For the GNN and more specifically GCNs, you can look at the original blog post by the author of the GCN paper: <https://tkipf.github.io/graph-convolutional-networks/>.

Here is also a short youtube video describing and visualizing the basics of GCN <https://youtu.be/2KRAOZIULzw>. Note that his aggregation and update of the layer is a bit different. But the visualization is great for understanding the concept of Graph Neural Networks.

1 Parsing Context-Free Grammars

In this exercise we cover the classical CKY (Cocke-Younger-Kasami, sometimes also CYK) algorithm for parsing context-free grammars (CFG). The CFG is discussed in chapter 12.2 of the book, and the CKY algorithm is discussed in chapter 13.2 of the book, which can be used as a reference.

1.1 Question 1.A

We consider the following CFG:

$$\begin{aligned}
 S &\rightarrow NP \quad VP \\
 NP &\rightarrow D \quad N \quad | \quad NP \quad PP \quad | \quad Pro \\
 VP &\rightarrow V \quad NP \quad | \quad V \quad NP \quad PP \\
 PP &\rightarrow P \quad NP \\
 D &\rightarrow the \quad | \quad a \\
 P &\rightarrow with \\
 Pro &\rightarrow john \\
 N &\rightarrow cat \quad | \quad tree \quad | \quad binoculars \\
 V &\rightarrow sees \quad | \quad climbs
 \end{aligned}$$

Now answer the solve the following open questions:

1. Think of two sentences that would be grammatical according to this grammar.

1.2 Question 1.B

We would like to use this grammar to parse sentences using the CKY algorithm. However, for CKY the grammar needs to be in Chomsky Normal Form (CNF). Convert the grammar given above into CNF.

1.3 Question 1.C

Now that we have our grammar in CNF, we could parse sentences using CKY. Find **all the trees** for the following sentences:

1. John climbs the tree.
2. The cat sees the tree with the binoculars.

1.4 Question 1.D

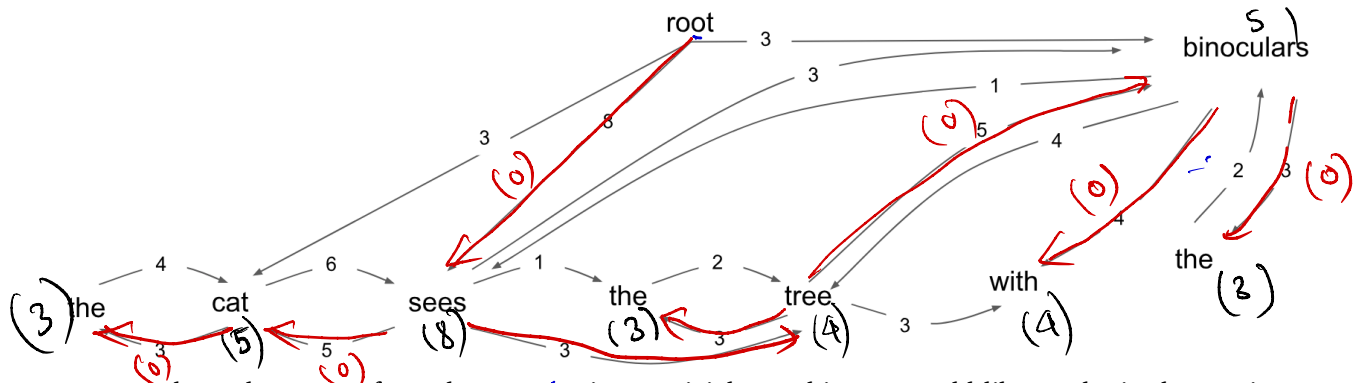
If everything went correct, you found two possible trees for sentence 2.

This type of ambiguity is called PP-attachment ambiguity. Make use of the difference in trees to describe what this ambiguity is and when it occurs.

2 Dependency Parsing

In this section we cover a classical algorithm for edge-factored dependency parsing models. This topic is covered in Chapter 14 of the book (particularly 14.5).

We assume we already have an edge scoring model. Consider the following scores given by our model:



To get a dependency tree from these scores is not trivial. For this, we would like to obtain the maximum spanning tree based on these scores. Use the Chu-Liu Edmonds (CLE) algorithm to get the maximum spanning tree. To which of the two meanings does this dependency parse correspond?

3 Using the dependency tree in a Graph Neural Network

In this exercise we work with a simplified model from the "Graph Convolution over Pruned Dependency Trees Improves Relation Extraction" paper: <https://aclanthology.org/D18-1244/>.

We just extracted a dependency tree in Question 2 that allows us to get insight into the meaning of the sentence. In this question, we will see how we can make use of the dependency tree to make other predictions for the sentence.

We are going to do semantic role labelling (SRL). We define a very simplistic model. The definition of the model is:

$$\hat{y} = \text{softmax}(F_{\text{pred}}([h; p]))$$

Where $;$ means concatenation, $h \in \mathbb{R}^e$ is the vector representation of a word, $\hat{y} \in \mathbb{R}^l$ the predicted label probabilities, and p the boolean value indicating if the word is the predicate or not, and F_{pred} is the prediction function or model. Our model consists of a GCN and as final step we use a linear layer function with the weight matrix $W_{\text{pred}} \in \mathbb{R}^{m \times l}$ and biases $b_{\text{pred}} \in \mathbb{R}^m$. e is the embedding size, l is the number of labels, m is the number hidden dimension. The possible labels are a reduced version of the prob bank: <ARG0, ARG1, O> with ARG0 the agent, ARG1 the patient, and O not part of a role.

There are many options for h , such as pretrained word embeddings (word2vec, BERT) or some model like a bi-LSTM. However, we want stronger features that indicate relations between words. So we use a GCN over the dependency relations, to update the hidden states before passing them through the prediction layer. In addition to the dependency relations, we also use bidirectional sequence relations: each word has an additional incoming edge from the previous word and from the following word.

We don't have edge labels, so we don't use them. We define the direction of the edge as two different labels (incoming/outgoing), so the model can learn a distinction between incoming dependencies and outgoing dependencies.

Here is the definition of the GCN (taken and adapted from Equation 1 from the paper mentioned above):

$$h_i^{(t+1)} = \text{ReLU}(h^{(t)} + \sum_{n \in \mathbb{N}(i)} (W_{in}^T h_n + b_{in}) + \sum_{m \in \mathbb{M}(i)} (W_{out}^T h_m + b_{out}))$$

With \mathbb{N} and \mathbb{M} the set of incoming and outgoing neighbouring nodes respectively. $W_{in} \in \mathbb{R}^{m \times m}$ and $W_{out} \in \mathbb{R}^{m \times m}$ the weight matrices for incoming and outgoing neighbours, and $b_{in} \in \mathbb{R}^m$ and $b_{out} \in \mathbb{R}^m$ the

biases for incoming and outgoing neighbours. m is the hidden size.

We chose the hidden size between timesteps or GCN layers to be the same. This makes it possible to reuse the same weight matrices for multiple layers. It is not mandatory to keep the hidden size equal. In fact, for convolutions it is common to reduce size after every layer. In literature you also see different weight matrices for every layer. In this question we simply use a single timestep (or single GCN layer).

NOTE: In this exercise we have multiple types of connections, without distinction between them (so we don't use their labels and don't apply different weights). This can result in multiple identical connections between words.

NOTE2: We did not use the default GCN computations with adjacency matrices and such, but we used the version that takes the sum over the neighbours. This simplifies the understanding and makes it interpretable as a message passing algorithm.

Question: Predict the labels of each word using SRL model with GCN for the sentence:

the cat sees the tree with the binoculars

As dependency tree, use the prediction from Question 2. The predicate of this sentence is **sees**.

Use the given weights and biases below. We also provide the initial embeddings for the words in the sentence in matrix $E \in \mathbb{R}^{|V| \times 2}$, where $|V|$ is the vocabulary size. The vocabulary is:

['binoculars', 'cat', 'sees', 'the', 'tree', 'with']

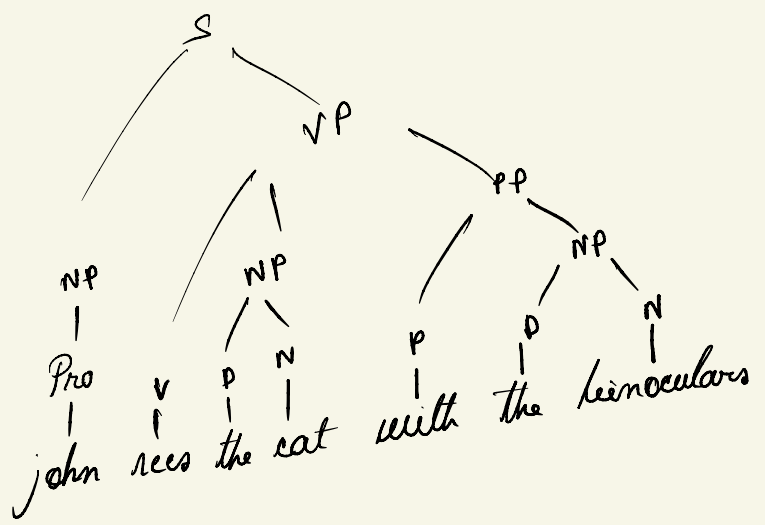
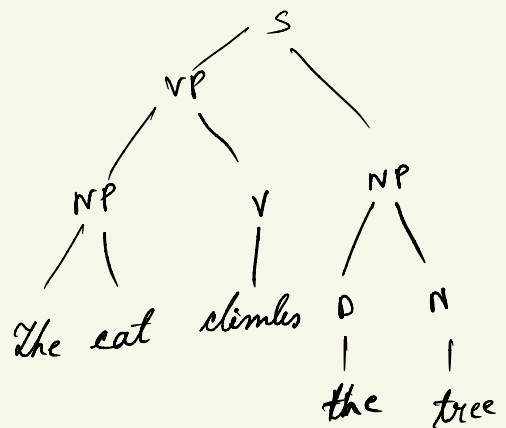
Use the following steps to solve the question:

1. For each node, define the two sets of neighbours, incoming and outgoing. Duplicates are allowed.
2. Do one GCN step and update the hidden state for all the nodes in the graph to $h^{(1)}$. Use the word embedding as $h^{(0)}$, for each node. Do not forget to append the predicate boolean p to the hidden state before processing it through the GCN.
3. Use $h^{(1)}$ to predict the role label for each of the words in the sentence.

$$E = \begin{bmatrix} -0.5 & 0.7 \\ 1.4 & 1.1 \\ 0.8 & -0.5 \\ -0.3 & 0.1 \\ 1.3 & 0.8 \\ 0.8 & -0.8 \end{bmatrix} \quad b_{in} = \begin{bmatrix} 0.6 \\ 3.0 \\ -1.0 \end{bmatrix} \quad b_{pred} = \begin{bmatrix} 1.0 \\ -0.5 \\ -0.5 \end{bmatrix} \quad b_{out} = \begin{bmatrix} -0.7 \\ 1.0 \\ 0.2 \end{bmatrix}$$

$$W_{in} = \begin{bmatrix} -0.5 & 0.0 & 0.6 \\ -2.0 & 0.5 & 1.0 \\ 2.0 & -0.7 & -0.8 \end{bmatrix} \quad W_{pred} = \begin{bmatrix} 3.0 & -2.0 & 0.2 \\ 0.1 & 2.0 & -3.0 \\ 0.7 & -2.0 & 2.0 \end{bmatrix} \quad W_{out} = \begin{bmatrix} 0.1 & -1.0 & 1.0 \\ -0.7 & -2.0 & 2.0 \\ 1.0 & 0.3 & -0.8 \end{bmatrix}$$

①

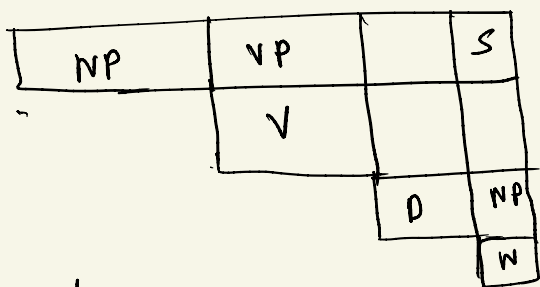


1.1

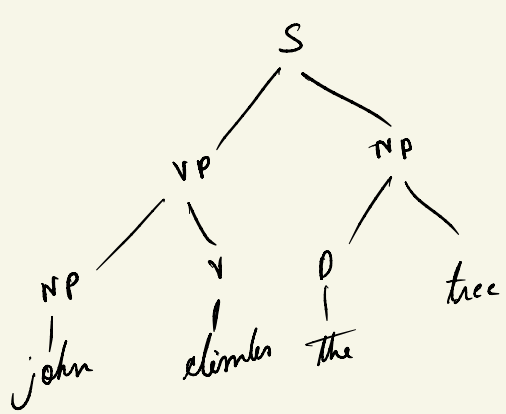
$S \rightarrow NP VP$
 $NP \rightarrow D N \mid NP PP \mid \text{john}$
 $VP \rightarrow V NP \mid X PP$
 $PP \rightarrow P NP \mid X \rightarrow V NP$

$D \rightarrow \text{the} \mid a$
 $P \rightarrow \text{with}$

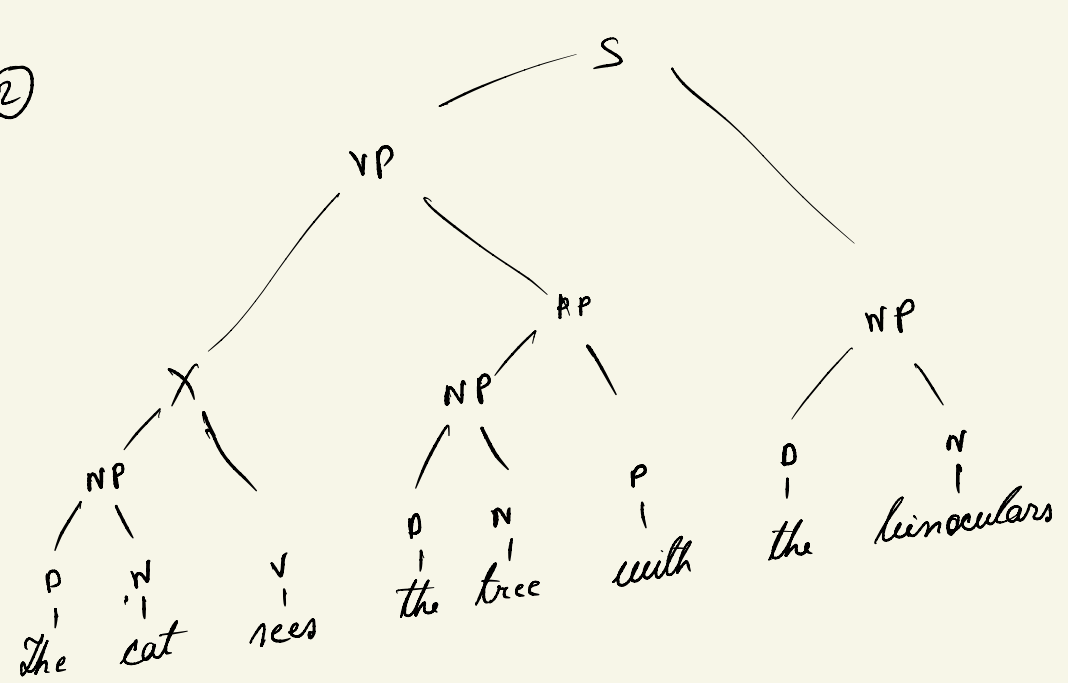
$N \rightarrow \text{cat} \mid \text{tree} \mid \text{binoculars}$
 $V \rightarrow \text{sees} \mid \text{climbs}$



John climbs the tree



②



① $G=(V,E)$

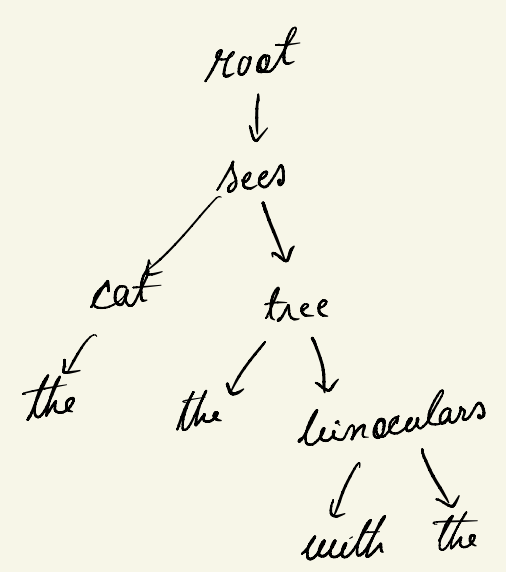
$X \in \mathbb{R}^{N \times D}, A \in \mathbb{R}^{N \times N}$

$A \times W$

$$\therefore f(H^{(k)}, A) = O(A H^{(k)} W^{(k)})$$

$$\hat{y} = \text{softmax}(F_{\text{pred}}([h; P]))$$

$$h_i^{t+1} = \text{RELU}\left(h^{(t)} + \sum_{n \in N(i)} (W_{in}^T h_n + b_{in}) + \sum_{m \in M(i)} (W_{out}^T h_m + b_{out})\right)$$



	binoculars	cat	sees	the	tree	with
binoculars	0	0	0	1	1	1
cat	0	0	1	1	0	0
sees	0	1	0	0	1	0
the	1	1	0	1	0	0
tree	1	0	1	1	0	0
with	1	0	0	0	0	0