

1 Preliminaries: Constituent and Dependency Trees

Let $w = (w_1, \dots, w_L)$ be a sentence.

A *constituent tree* is a rooted tree whose leaves are the words $(w_i)_{i=1}^N$ and internal nodes are constituents.

A *constituent* is a triple (Z, \mathcal{Y}, h) containing, respectively, its label, yield, and lexical head which satisfy some constraints¹.

A constituent is *discontinuous* if its yield is not contiguous.

A *dependency tree* is a rooted tree spanning the words in the sentence $(w_i)_{i=1}^N$. Each edge is labelled and connects a parent word (head) to a child word (dependency).

Fernández-González and Martins (2015) show that under quite general conditions², constituent trees are isomorphic to dependency trees in which the edges contain information about constituent labels and attachment order.

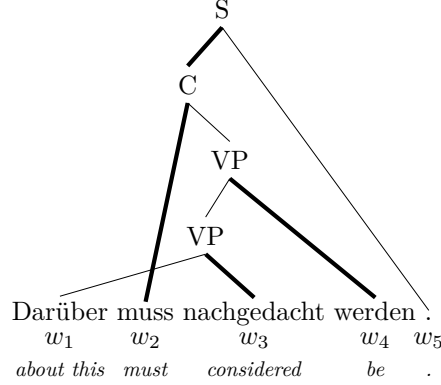


Figure 1.1: Example of a discontinuous constituent tree for the German sentence ‘*Darüber muss nachgedacht werden.*’ (‘this must be considered.’). Bold lines indicate head words.

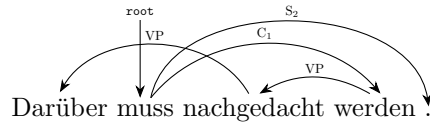
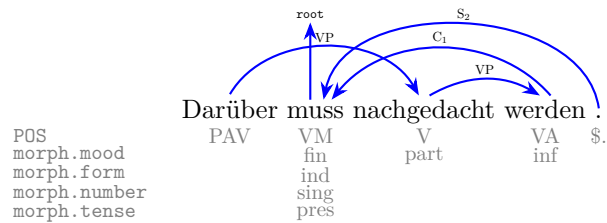


Figure 1.2: The dependency tree corresponding to the constituent tree in Fig. 1.1.

2 Mathematical Model



Given regressor $w = (w_i)_{i=1}^N$, our goal is to model its dependency tree along with information like its part-of-speech (POS) and morphology.

Our parser will work from the bottom-up, so we will think of *arcs* going from every child w_i to its parent w_i^{arc} .

Denote the regressand by $y = (y_i)_{i=1}^N$ where $y_i = (w_i^{\text{arc}}, w_i^{\text{lab}}, w_i^{\text{ord}}, w_i^{\text{pos}}, w_i^{\text{morph}})$. Define $y_{<i} = (y_1, \dots, y_i)$ for each $i = 2, \dots, N$ and $y_{<1} = 0$.

¹A leaf node can be modelled as a constituent whose yield contains only its head. For every constituent (Z, \mathcal{Y}, h) with children $\{(A_k, \mathcal{X}_k, m_k)\}$, 1. $\mathcal{Y} = \bigcup_{k=1}^K \mathcal{X}_k$, and 2. there is a unique k such that $h = m_k$.

²The constituent trees must be unaryless.

Assumption. For each $i = 1, \dots, N$, the random variables $w_i^{\text{lab}}, w_i^{\text{ord}}, w_i^{\text{pos}}$ and w_i^{morph} are mutually independent conditional on w_i^{arc}, w and $y_{<i}$.³

We can decompose the conditional probability of y given w :

$$p(y | w) = \prod_{i=1}^N p(y_i | y_{<i}, w) \quad (2.1)$$

$$= \prod_{i=1}^N p(w_i^{\text{arc}} | y_{<i}, w) p(w_i^{\text{lab}} | w_i^{\text{arc}}, y_{<i}, w) p(w_i^{\text{ord}} | w_i^{\text{arc}}, y_{<i}, w) p(w_i^{\text{pos}} | w_i^{\text{arc}}, y_{<i}, w) p(w_i^{\text{morph}} | w_i^{\text{arc}}, y_{<i}, w). \quad (2.2)$$

3 Encoder-Decoder Setup

Let $N = \{1, \dots, n\}$ be an index set. Given an input sentence $w = (w_i)_{i=1}^N$ we generate a sequence of *embeddings* $\omega = (\omega_i)_{i=1}^N$ where

$$\omega_i = \text{WordEmbed}(w_i) \oplus \text{CharEmbed}(w_i) \oplus \text{BertEmbed}(w_i).$$

Character-level embeddings are implemented via a CNN á la Chiu and Nichols (2016). BERT embeddings are finetuned from a BERT model pre-trained on German text by Chan et al. (2020).

Encoder: feed embeddings through a multi-layer bi-directional LSTM with skip-connections and dropout:

$$\mathbf{e} = (\mathbf{e}_i)_{i=0, \dots, n} = \text{BiLSTM}(\omega)$$

Decoder: feed embeddings through a single-layer uni-directional LSTM with dropout:

$$\mathbf{d} = (\mathbf{d}_i)_{i=1, \dots, n} = \text{LSTM}(\omega)$$

4 Bi-affine Attention Mechanism

We feed \mathbf{e} and \mathbf{d} through MLPs to produce sequences $(\mathbf{e}^{\text{arc}}, \mathbf{d}^{\text{arc}})$ of dimension-reduced vectors. These are fed into a bi-affine layer which produces latent features \mathbf{v}^{arc} that are then fed into an attention layer, resulting in logits corresponding to strength of an arc.

$$\mathbf{e}^{\text{arc}} = \text{MLP}_{\text{enc}}^{\text{arc}}(\mathbf{e}); \quad \mathbf{d}^{\text{arc}} = \text{MLP}_{\text{dec}}^{\text{arc}}(\mathbf{d}) \quad (4.1)$$

$$\mathbf{v}_{i,j}^{\text{arc}} = \text{BiAff}^{\text{arc}}(\mathbf{e}_i^{\text{arc}}, \mathbf{d}_j^{\text{arc}}) \quad (4.2)$$

$$:= \mathbf{e}_i^{\text{arcT}} \mathbf{U}_{\text{h-d}}^{\text{arc}} \mathbf{d}_j^{\text{arc}} + \mathbf{e}_i^{\text{arcT}} \mathbf{U}_{\text{h-h}}^{\text{arc}} \mathbf{e}_i^{\text{arc}} + \mathbf{d}_j^{\text{arcT}} \mathbf{U}_{\text{d-d}}^{\text{arc}} \mathbf{d}_j^{\text{arc}} \quad (4.3)$$

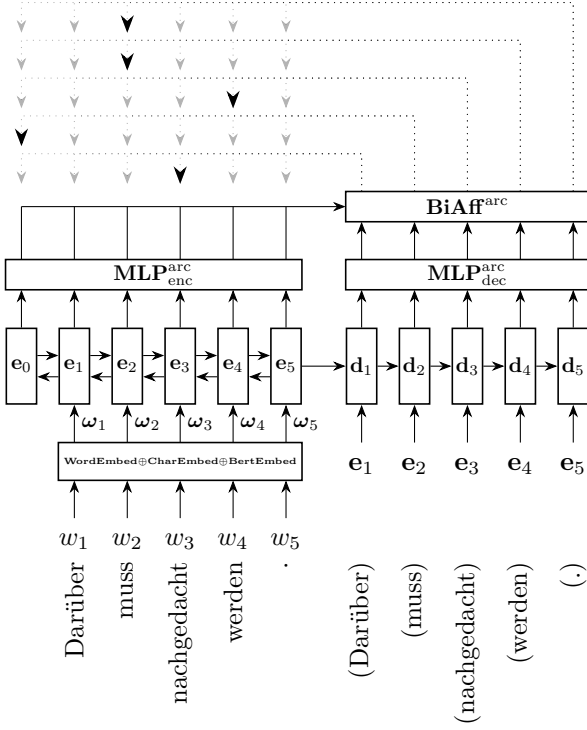
$$+ U_{\text{h}}^{\text{arc}} \mathbf{e}_i^{\text{arc}} + U_{\text{d}}^{\text{arc}} \mathbf{d}_j^{\text{arc}} + \mathbf{u}_{\text{bias}}^{\text{arc}} \quad (4.4)$$

$$s_{i,j}^{\text{arc}} = \mathbf{u}_{\text{agg}}^{\text{arcT}} \tanh(\mathbf{v}_{i,j}^{\text{arc}}) \quad (4.5)$$

Fixing dependency j , the vector $\text{softmax}(\mathbf{s}_{:,j})$ can be interpreted as an estimated probability distribution over potential heads.

$$\hat{p}^{\text{arc}}(w_i | y_{<j}, w) = \text{softmax}(\mathbf{s}_{:,j}^{\text{arc}})_i.$$

³If w_i^{morph} is a vector, assume every component is mutually independent with each other and with $w_i^{\text{lab}}, w_i^{\text{ord}}, w_i^{\text{pos}}$ conditional on $w_i^{\text{arc}}, y_{<i}, w$



5 Bi-affine Classifier for Attachment Order, POS and Morphology

Attachment order, POS and morphologies are predicted via a classification layer. We use a bi-affine classifier which allows us to model probabilities of classes *conditional* on arcs, and thus use structural cues in addition to encoder/decoder states to better capture the complexity of the language. The encoder and decoder are *shared* across the tasks.

For example suppose we would like to predict the part of speech $c \in \mathcal{C}$ for word w_j conditional on its parent being w_i .

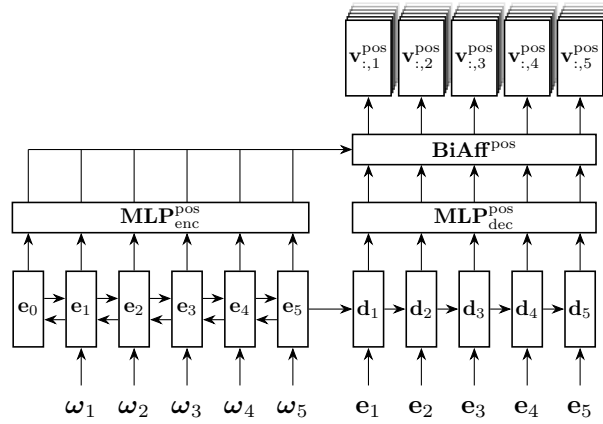
$$\mathbf{e}^{\text{pos}} = \text{MLP}_{\text{enc}}^{\text{pos}}(\mathbf{e}); \quad \mathbf{d}^{\text{pos}} = \text{MLP}_{\text{dec}}^{\text{pos}}(\mathbf{d}) \quad (5.1)$$

$$\mathbf{v}_{i,j}^{\text{pos}} = \text{BiAff}^{\text{pos}}(\mathbf{e}_i^{\text{pos}}, \mathbf{d}_j^{\text{pos}}) \quad (5.2)$$

$$:= \mathbf{e}_i^{\text{pos}\top} \mathbf{U}_{\text{h-d}}^{\text{pos}} \mathbf{d}_j^{\text{pos}} + \mathbf{e}_i^{\text{pos}\top} \mathbf{U}_{\text{h-h}}^{\text{pos}} \mathbf{e}_i^{\text{pos}} + \mathbf{d}_j^{\text{pos}\top} \mathbf{U}_{\text{d-d}}^{\text{pos}} \mathbf{d}_j^{\text{pos}} \quad (5.3)$$

$$+ U_{\text{h}}^{\text{pos}} \mathbf{e}_i^{\text{pos}} + U_{\text{d}}^{\text{pos}} \mathbf{d}_j^{\text{pos}} + \mathbf{u}_{\text{bias}}^{\text{pos}} \quad (5.4)$$

$$\hat{p}^{\text{pos}}(c \mid w_i; y_{<j}, w) = \text{softmax}(\mathbf{v}_{i,j}^{\text{pos}})_c \quad (5.5)$$



Bibliography

- Chan, B., Möller, T., Pietsch, M., & Soni, T. (2020). *German BERT (bert-base-german-cased)*. Retrieved September 1, 2023, from <https://huggingface.co/bert-base-german-cased>
- Chiu, J. P., & Nichols, E. (2016). Named entity recognition with bidirectional lstm-cnns. *Transactions of the association for computational linguistics*, 4, 357–370.
- Fernández-González, D., & Martins, A. F. (2015). Parsing as reduction. *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 1523–1533.