

# CS 224n Assignment #3: Dependency Parsing

Ming Liu

*Updated May, 11 at 11:10am*

In this assignment, you will build a neural dependency parser using PyTorch. For a review of the fundamentals of PyTorch, please check out the PyTorch review session on Canvas. You will implement and train a dependency parser, before analyzing a few erroneous dependency parses.

## 1. Neural Transition-Based Dependency Parsing (44 points)

In this section, you'll be implementing a neural-network based dependency parser with the goal of maximizing performance on the UAS (Unlabeled Attachment Score) metric.

Before you begin, please follow the README to install all the needed dependencies for the assignment. We will be using PyTorch 1.7.1 from <https://pytorch.org/get-started/locally/> with the CUDA option set to None, and the tqdm package – which produces progress bar visualizations throughout your training process. The official PyTorch website is a great resource that includes tutorials for understanding PyTorch's Tensor library and neural networks.

A dependency parser analyzes the grammatical structure of a sentence, establishing relationships between *head* words, and words which modify those heads. There are multiple types of dependency parsers, including transition-based parsers, graph-based parsers, and feature-based parsers. Your implementation will be a *transition-based* parser, which incrementally builds up a parse one step at a time. At every step it maintains a *partial parse*, which is represented as follows:

- A stack of words that are currently being processed.
- A buffer of words yet to be processed.
- A list of *dependencies* predicted by the parser.

Initially, the stack only contains ROOT, the dependencies list is empty, and the buffer contains all words of the sentence in order. At each step, the parser applies a *transition* to the partial parse until its buffer is empty and the stack size is 1. The following transitions can be applied:

- **SHIFT**: removes the first word from the buffer and pushes it onto the stack.
- **LEFT-ARC**: marks the second (second most recently added) item on the stack as a dependent of the first item and removes the second item from the stack, adding a *first\_word* → *second\_word* dependency to the dependency list.
- **RIGHT-ARC**: marks the first (most recently added) item on the stack as a dependent of the second item and removes the first item from the stack, adding a *second\_word* → *first\_word* dependency to the dependency list.

On each step, your parser will decide among the three transitions using a neural network classifier.

- (a) (4 points) Complete the sequence of transitions needed for parsing the sentence “*Today I parsed a sentence*”. The dependency tree for the sentence is shown below. At each step, give the configuration of the stack and buffer, as well as what transition was applied this step and what new dependency was added (if any). The first four steps are provided below as an example.



Stack	Buffer	New dependency	Transition
[ROOT]	[Today, I, parsed, a, sentence]		Initial Configuration
[ROOT, Today]	[I, parsed, a, sentence]		SHIFT
[ROOT, Today, I]	[parsed, a, sentence]		SHIFT
[ROOT, Today, I, parsed]	[a, sentence]		SHIFT
[ROOT, Today, parsed]	[a, sentence]	parsed→I	LEFT-ARC

**Solution:**

Here are the configurations in following steps of the sentence's parsing.

Stack	Buffer	New dependency	Transition
[ROOT, Today, parsed, a]	[sentence]		SHIFT
[ROOT, Today, parsed, a, sentence]	[]		SHIFT
[ROOT, Today, parsed, sentence]	[]	sentence→a	LEFT-ARC
[ROOT, Today, parsed]	[]	parsed→sentence	RIGHT-ARC
[ROOT, parsed]	[]	parsed→Today	LEFT-ARC
[ROOT]	[]	ROOT→parsed	LEFT-ARC
[ROOT]	[]		DONE

- (b) (2 points) A sentence containing  $n$  words will be parsed in how many steps (in terms of  $n$ )? Briefly explain in 1-2 sentences why.

**Solution:**

In *transition-based* parsing, the stack only contains ROOT and the buffer is empty when the parser finished parsing a sentence. Accordingly, we could judge whether a sentence is parsed by our parser by this condition:  $len(stack) == 1$  and  $len(buffer) == 0$ .

Therefore, a sentence containing  $n$  words will be parsed in  $O(n)$  steps.

- (c) (6 points) Implement the `__init__` and `parse_step` functions in the `PartialParse` class in `parser_transitions.py`. This implements the transition mechanics your parser will use. You can run basic (non-exhaustive) tests by running `python parser_transitions.py part.c`.

**Solution:**

Here are the implementation of `__init__` and `parse_step` functions in the `PartialParse` class.

```
def __init__(self, sentence):
    """Initializes this partial parse.
    """
    self.sentence = sentence
    self.stack = ["ROOT"]
    self.buffer = self.sentence.copy()
    self.dependencies = []
```

Listing 1: Implementation of `__init__` function in `PartialParse` class.

```

def parse_step(self, transition):
    """Performs a single parse step by applying the given
    transition to this partial parse.
    """
    if transition == "S":
        self.stack.append(self.buffer[0])
        self.buffer = self.buffer[1:]
    elif transition == "LA":
        self.dependencies.append((self.stack[-1], self.stack[-2]))
        self.stack.pop(-2)
    elif transition == "RA":
        self.dependencies.append((self.stack[-2], self.stack[-1]))
        self.stack.pop()
    else:
        raise Exception("Unknown transition %s" % transition)

```

Listing 2: Implementation of parse\_step function in PartialParse class.

- (d) (8 points) Our network will predict which transition should be applied next to a partial parse. We could use it to parse a single sentence by applying predicted transitions until the parse is complete. However, neural networks run much more efficiently when making predictions about *batches* of data at a time (i.e., predicting the next transition for many different partial parses simultaneously). We can parse sentences in minibatches with the following algorithm.

---

#### Algorithm 1 Minibatch Dependency Parsing

---

**Input:** sentences, a list of sentences to be parsed and model, our model that makes parse decisions

Initialize partial\_pares as a list of PartialPares, one for each sentence in sentences

Initialize unfinished\_pares as a shallow copy of partial\_pares

**while** unfinished\_pares is not empty **do**

    Take the first batch\_size parses in unfinished\_pares as a minibatch

    Use the model to predict the next transition for each partial parse in the minibatch

    Perform a parse step on each partial parse in the minibatch with its predicted transition

    Remove the completed (empty buffer and stack of size 1) parses from unfinished\_pares

**end while**

**Return:** The dependencies for each (now completed) parse in partial\_pares.

---

Implement this algorithm in the minibatch\_parse function in parser\_transitions.py. You can run basic (non-exhaustive) tests by running python parser\_transitions.py part\_d.

*Note: You will need minibatch\_parse to be correctly implemented to evaluate the model you will build in part (e). However, you do not need it to train the model, so you should be able to complete most of part (e) even if minibatch\_parse is not implemented yet.*

**Solution:**

Here are the implementation of algorithm 1 by which we could parse sentences in minibatches.

```
def minibatch_parse(sentences, model, batch_size):
    """Parses a list of sentences in minibatches using a model.
    """
    dependencies = []

    partial_parses = [PartialParse(sent) for sent in sentences]
    # unfinished_parses = partial_parses.copy()
    unfinished_parses = partial_parses[:]

    while len(unfinished_parses) > 0:
        minibatch = unfinished_parses[:batch_size].copy()
        transitions = model.predict(unfinished_parses[:batch_size])

        for i, batch_parse in enumerate(minibatch):
            batch_parse.parse_step(transitions[i])
            # remove the parsed ones
            if len(batch_parse.stack) == 1 and \
               len(batch_parse.buffer) == 0:
                unfinished_parses.remove(batch_parse)

    dependencies = [parse.dependencies for parse in partial_parses]

    return dependencies
```

Listing 3: Implementation of minibatch\_parse function in parser\_transitions.py.

- (e) (12 points) We are now going to train a neural network to predict, given the state of the stack, buffer, and dependencies, which transition should be applied next.

First, the model extracts a feature vector representing the current state. We will be using the feature set presented in the original neural dependency parsing paper: *A Fast and Accurate Dependency Parser using Neural Networks*.<sup>1</sup> The function extracting these features has been implemented for you in `utils/parser_utils.py`. This feature vector consists of a list of tokens (e.g., the last word in the stack, first word in the buffer, dependent of the second-to-last word in the stack if there is one, etc.). They can be represented as a list of integers  $\mathbf{w} = [w_1, w_2, \dots, w_m]$  where  $m$  is the number of features and each  $0 \leq w_i < |V|$  is the index of a token in the vocabulary ( $|V|$  is the vocabulary size). Then our network looks up an embedding for each word and concatenates them into a single input vector:

$$\mathbf{x} = [\mathbf{E}_{w_1}, \dots, \mathbf{E}_{w_m}] \in \mathbb{R}^{dm}$$

where  $\mathbf{E} \in \mathbb{R}^{|V| \times d}$  is an embedding matrix with each row  $\mathbf{E}_w$  as the vector for a particular word  $w$ . We then compute our prediction as:

$$\begin{aligned} \mathbf{h} &= \text{ReLU}(\mathbf{XW} + \mathbf{b}_1) \\ \mathbf{l} &= \mathbf{hU} + \mathbf{b}_2 \\ \hat{\mathbf{y}} &= \text{softmax}(\mathbf{l}) \end{aligned}$$

<sup>1</sup>Chen and Manning, 2014, <https://nlp.stanford.edu/pubs/emnlp2014-depparser.pdf>

$\mathbf{X}$  is a mini-batch of embedded inputs of shape (batch\_size, dm).  $\mathbf{h}$  is the hidden layer activation of shape (batch\_size, hidden\_size).  $\mathbf{W}$  and  $\mathbf{b}_1$  are the weight matrix and bias vector which transform  $\mathbf{x}$  into  $\mathbf{h}$ . And  $\text{ReLU}(z) = \max(z, 0)$ .

$\mathbf{l}$  is the matrix of output logits in shape (batch\_size, num\_classes).  $\mathbf{U}$  and  $\mathbf{b}_2$  are the weight matrix and bias vector which transform  $\mathbf{h}$  into  $\mathbf{l}$ .

Finally,  $\hat{\mathbf{y}}$  is the model's final prediction in shape (batch\_size, num\_classes). Each row is a probability distribution (sums up to 1) over all classes.

We will train the model to minimize cross-entropy loss:

$$J(\theta) = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^3 y_i \log \hat{y}_i$$

To compute the loss for the training set, we average this  $J(\theta)$  across all training examples.

We will use UAS score as our evaluation metric. UAS stands for Unlabeled Attachment Score, which is computed as the ratio between the number of correctly predicted dependencies and the number of total dependencies. UAS score is “Unlabeled” because it ignores the types of the dependency relations, which our model does not predict.

In `parser_model.py` you will find skeleton code to implement this simple neural network using PyTorch. Complete the `__init__`, `embedding_lookup` and `forward` functions to implement the model. Then complete the `train_for_epoch` and `train` functions within the `run.py` file.

Finally execute `python run.py` to train your model and compute predictions on test data from Penn Treebank (annotated with Universal Dependencies).

**Solution:**

Here are the implementation of the functions to be completed in `parser_model.py` and `run.py` file.

```
def __init__(self, embeddings, n_features=36,
             hidden_size=200, n_classes=3, dropout_prob=0.5):
    super(ParserModel, self).__init__()
    self.n_features = n_features
    self.n_classes = n_classes
    self.dropout_prob = dropout_prob
    self.embed_size = embeddings.shape[1] # shape[0] is num_words
    self.hidden_size = hidden_size
    self.embeddings = nn.Parameter(torch.tensor(embeddings))

    self.embed_to_hidden_weight = nn.Parameter(
        nn.init.xavier_uniform_(torch.empty(
            self.n_features * self.embed_size, self.hidden_size
        )))
    self.embed_to_hidden_bias = nn.Parameter(nn.init.uniform_(
        torch.empty(self.hidden_size,)
    ))
    self.dropout = nn.Dropout(p=self.dropout_prob)
    self.hidden_to_logits_weight = nn.Parameter(
        nn.init.xavier_uniform_(
            torch.empty(self.hidden_size, self.n_classes)
        ))
    self.hidden_to_logits_bias = nn.Parameter(nn.init.uniform_(
        torch.empty(self.n_classes,)
    ))
```

Listing 4: Implementation of `__init__` function in `parser_model.py`.

```
def embedding_lookup(self, w):
    """ Utilize `w` to select embeddings from
        embedding matrix `self.embeddings`

        @param w (Tensor): input tensor of word indices
                           (batch_size, n_features)

        @return x (Tensor): tensor of embeddings for words represented
                           in w (batch_size, n_features * embed_size)
    """

    # embeddings: word embeddings (num_words, embedding_size)
    x = torch.index_select(self.embeddings, 0, w.view(-1))
    x = x.view(w.shape[0], -1)

    return x
```

Listing 5: Implementation of `embedding_lookup` function in `parser_model.py`.

```

def forward(self, w):
    """ Run the model forward.
    """

    embedding_res = self.embedding_lookup(w)
    affine_res = torch.matmul(embedding_res, \
                              self.embed_to_hidden_weight) \
                + self.embed_to_hidden_bias

    h = self.dropout(F.relu(affine_res))
    logits = torch.matmul(h, self.hidden_to_logits_weight) \
            + self.hidden_to_logits_bias

    return logits

```

Listing 6: Implementation of forward function in parser\_model.py.

```

def train_for_epoch(parser, train_data, dev_data,
                    optimizer, loss_func, batch_size):
    parser.model.train()
    n_minibatches = math.ceil(len(train_data) / batch_size)
    loss_meter = AverageMeter()

    with tqdm(total=(n_minibatches)) as prog:
        for i, (train_x, train_y) in \
            enumerate(minibatches(train_data, batch_size)):
            optimizer.zero_grad()
            loss = 0.
            train_x = torch.from_numpy(train_x).long()
            train_y = torch.from_numpy(train_y.nonzero()[1]).long()

            logits = parser.model(train_x)
            loss = loss_func(logits, train_y)
            loss.backward()
            optimizer.step()

            prog.update(1)
            loss_meter.update(loss.item())

    print ("Average Train Loss: {}".format(loss_meter.avg))

    print("Evaluating on dev set",)
    parser.model.eval()
    dev_UAS, _ = parser.parse(dev_data)
    print("- dev UAS: {:.2f}".format(dev_UAS * 100.0))
    return dev_UAS

```

Listing 7: Implementation of train\_for\_epoch function in run.py.

```

def train(parser, train_data, dev_data, output_path,
          batch_size=1024, n_epochs=10, lr=0.0005):
    """ Train the neural dependency parser.
    """
    best_dev_UAS = 0

    optimizer = optim.Adam(params=parser.model.parameters(), lr=lr)
    loss_func = nn.CrossEntropyLoss()

    for epoch in range(n_epochs):
        print("Epoch {:} out of {:}".format(epoch + 1, n_epochs))
        dev_UAS = train_for_epoch(parser, train_data, dev_data, \
                                   optimizer, loss_func, batch_size)
        if dev_UAS > best_dev_UAS:
            best_dev_UAS = dev_UAS
            print("New best dev UAS! Saving model.")
            torch.save(parser.model.state_dict(), output_path)
        print("")

```

Listing 8: Implementation of train function in parser\_model.py.

### Important Notes:

- For this assignment, you are asked to implement Linear layer and Embedding layer. Please **DO NOT** use **torch.nn.Linear** or **torch.nn.Embedding** module in your code, otherwise you will receive deductions for this problem.
- Please follow the naming requirements in our TODO if there are any, e.g. if there are explicit requirements about variable names you have to follow them in order to receive full credits. You are free to declare other variable names if not explicitly required.

### More Hints:

#### Implementation details:

- Each of the variables you are asked to declare (`self.embed_to_hidden_weight`, `self.embed_to_hidden_bias`, `self.hidden_to_logits_weight`, `self.hidden_to_logits_bias`) corresponds to one of the variables above (**W**, **b<sub>1</sub>**, **U**, **b<sub>2</sub>**).
- It may help to work backwards in the algorithm (start from **ŷ**) and keep track of the matrix/vector sizes.

#### Debugging help:

- Once you have implemented `embedding_lookup` (e) or `forward` (f) you can call `python parser_model.py` with flag `-e` or `-f` or both to run sanity checks with each function. These sanity checks are fairly basic and passing them doesn't mean your code is bug free.
- When debugging, you can add a debug flag: `python run.py -d`. This will cause the code to run over a small subset of the data, so that training the model won't take as long. Make sure to remove the `-d` flag to run the full model once you are done debugging.

#### Sanity checks:

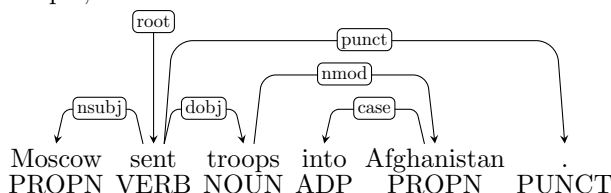
- In debug mode, you should be able to get a loss smaller than 0.2 and a UAS larger than 65 on the dev set (although in rare cases your results may be lower as there is some randomness when training).
- When debug mode is disabled, it should take about **1 hour** to train the model on the entire the training dataset.



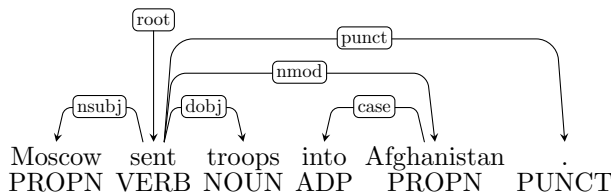
- When debug mode is disabled, you should be able to get a loss smaller than 0.08 on the train set and an Unlabeled Attachment Score larger than 87 on the dev set. For comparison, the model in the original neural dependency parsing paper gets 92.5 UAS. If you want, you can tweak the hyperparameters for your model (hidden layer size, hyperparameters for Adam, number of epochs, etc.) to improve the performance (but you are not required to do so).

### Deliverables:

- Working implementation of the transition mechanics that the neural dependency parser uses in `parser_transitions.py`.
  - Working implementation of minibatch dependency parsing in `parser_transitions.py`.
  - Working implementation of the neural dependency parser in `parser_model.py`. (We'll look at and run this code for grading).
  - Working implementation of the functions for training in `run.py`. (We'll look at and run this code for grading).
  - **Report the best UAS your model achieves on the dev set and the UAS it achieves on the test set in your writeup.**
- (f) (12 points) We'd like to look at example dependency parses and understand where parsers like ours might be wrong. For example, in this sentence:



the dependency of the phrase *into Afghanistan* is wrong, because the phrase should modify *sent* (as in *sent into Afghanistan*) not *troops* (because *troops into Afghanistan* doesn't make sense). Here is the correct parse:



More generally, here are four types of parsing error:

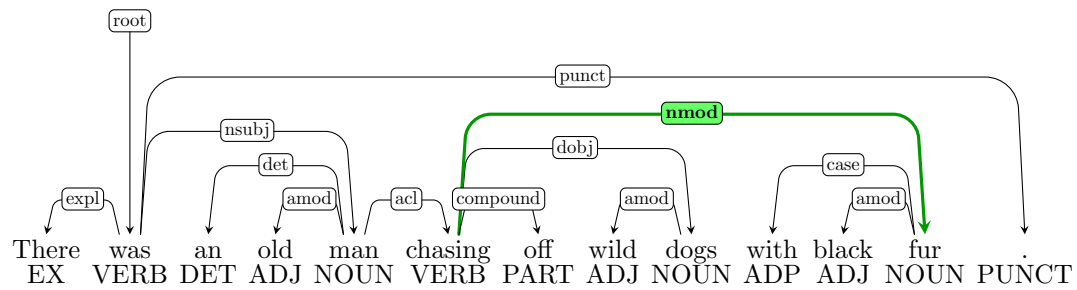
- **Prepositional Phrase Attachment Error:** In the example above, the phrase *into Afghanistan* is a prepositional phrase<sup>2</sup>. It modifies *sent*, specifying the destination of this action. Therefore, the correct dependency is *sent* → *Afghanistan*. A Prepositional Phrase Attachment Error is when a prepositional phrase is attached to the wrong head word. More examples of prepositional phrases include *with a rock*, *before midnight* and *under the carpet*.
- **Verb Phrase Attachment Error:** In the sentence *Leaving the store unattended, I went outside to watch the parade*, the phrase *leaving the store unattended* is a verb phrase<sup>3</sup>. In this example, this verb phrase modifies *went* (*went* → *leaving*). A Verb Phrase Attachment Error is when a verb phrase is attached to the wrong head word.
- **Modifier Attachment Error:** In the sentence *I am extremely short*, the adverb *extremely* is a modifier of the adjective *short*. The correct head word of *extremely* is *short* (*short* → *extremely*). A Modifier Attachment Error is when a modifier is attached to the wrong head word.

<sup>2</sup>For examples of prepositional phrases, see: <https://www.grammarly.com/blog/prepositional-phrase/>

<sup>3</sup>For examples of verb phrases, see: <https://examples.yourdictionary.com/verb-phrase-examples.html>

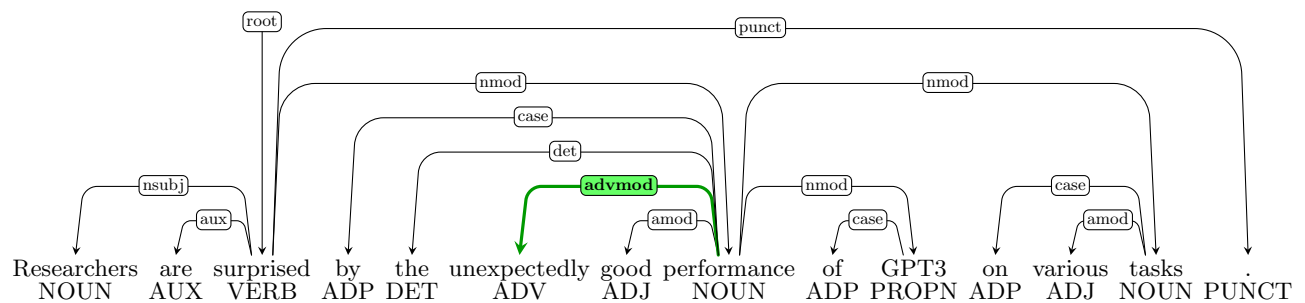


ii.

**Solution:**

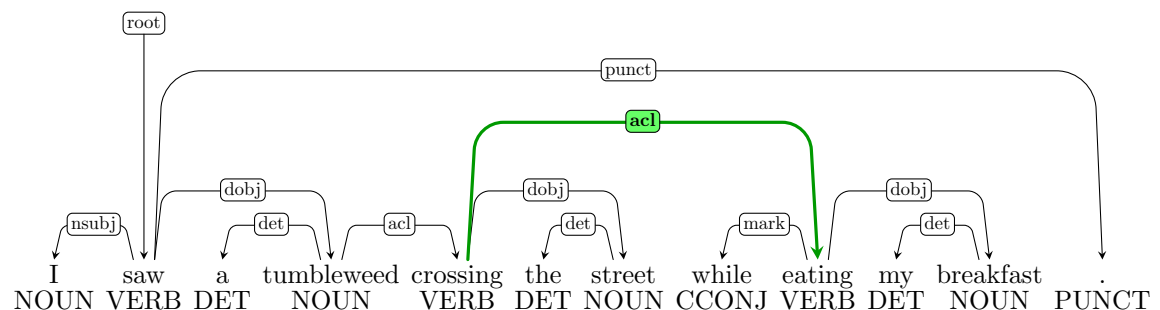
- **Error type:** Prepositional Phrase Attachment Error
- **Incorrect dependency:** chasing → fur
- **Correct dependency:** dogs → fur

iii.

**Solution:**

- **Error type:** Modifier Attachment Error
- **Incorrect dependency:** performance → unexpectedly
- **Correct dependency:** good → unexpectedly

iv.

**Solution:**

- **Error type:** Verb Phrase Attachment Error
- **Incorrect dependency:** crossing → eating
- **Correct dependency:** saw → eating