

# Neural language models with latent syntax

---

**Daan van Stigt**

Institute for Logic, Language and Computation

# Contents

---

<b>Acknowledgements</b>	<b>2</b>
<b>1 Introduction</b>	<b>4</b>
<b>2 Background</b>	<b>7</b>
2.1 Syntax . . . . .	7
2.2 Parsing . . . . .	12
2.3 Language models . . . . .	17
2.4 Neural networks . . . . .	20
<b>3 Recurrent Neural Network Grammars</b>	<b>22</b>
3.1 Model . . . . .	22
3.2 Parametrization . . . . .	27
3.3 Training . . . . .	29
3.4 Inference . . . . .	29
3.5 Experiments . . . . .	31
3.6 Related work . . . . .	33
<b>4 Conditional Random Field parser</b>	<b>34</b>
4.1 Model . . . . .	35
4.2 Parametrization . . . . .	36
4.3 Inference . . . . .	37

4.4	Experiments . . . . .	44
<b>5</b>	<b>Semisupervised learning</b>	<b>48</b>
5.1	Objective . . . . .	49
5.2	Variational approximation . . . . .	50
5.3	Optimization . . . . .	51
5.4	Variance reduction . . . . .	52
5.5	Experiments . . . . .	52
5.6	Related work . . . . .	53
<b>6</b>	<b>Syntactic evaluation</b>	<b>54</b>
6.1	Syntactic evaluation . . . . .	55
6.2	Multitask learning . . . . .	57
6.3	Experiments . . . . .	60
<b>7</b>	<b>Conclusion</b>	<b>62</b>
7.1	Main contributions . . . . .	62
7.2	Future work . . . . .	62
	<b>Appendices</b>	<b>63</b>
<b>A</b>	<b>Figures</b>	<b>64</b>
A.1	Training plots . . . . .	64
A.2	Test scores . . . . .	65
A.3	Samples . . . . .	65
A.4	Syntactic evaluation . . . . .	65
<b>B</b>	<b>Implementation</b>	<b>66</b>
B.1	Data . . . . .	66
B.2	Implementation . . . . .	69
<b>C</b>	<b>Semiring parsing</b>	<b>71</b>
C.1	Hypergraph . . . . .	71
C.2	Semiring . . . . .	72
C.3	Semiring parsing . . . . .	74

<b>D</b>	<b>Variational inference</b>	<b>78</b>
D.1	Variational Inference . . . . .	78
D.2	Score function estimator . . . . .	79
D.3	Optimization . . . . .	80
D.4	Variance reduction . . . . .	81
<b>E</b>	<b>Syneval dataset</b>	<b>86</b>
	<b>References</b>	<b>89</b>

---

## ABSTRACT

In this thesis I investigate the question: *What are effective ways of incorporating syntactic structure into neural language models?*

In this thesis I:

- study a class of neural language models that merges generative transition-based parsing with recurrent neural networks in order to model sentences together with their latent syntactic structure;
  - propose a new globally trained chart-based parser as an alternative proposal distribution used in the approximate marginalization;
  - propose effective methods for semisupervised learning, making the syntactic structure a latent variable;
  - perform targeted syntactic evaluation and compare the model's performance with that of alternative models that are based on multitask learning.
-

## Acknowledgements

---

$\mathbf{a}, \mathbf{b}, \dots$	Vectors over the reals, <i>i.e.</i> $\mathbf{a} \in \mathbb{R}^m$ .
$\mathbf{A}, \mathbf{B}, \dots$	Matrices over the reals, <i>i.e.</i> $\mathbf{A} \in \mathbb{R}^{m \times n}$ .
$[\mathbf{a}]_i$	Vector indexing: $[\mathbf{a}]_i \in \mathbb{R}$ for $1 \leq i \leq m$ .
$[\mathbf{a}; \mathbf{b}]$	Vector concatenation: $\mathbf{a} \in \mathbb{R}^m$ , $\mathbf{b} \in \mathbb{R}^n$ , $[\mathbf{a}; \mathbf{b}] \in \mathbb{R}^{m+n}$ .
$[\mathbf{a}, \mathbf{b}]$	Vertical vector stacking: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$ , $[\mathbf{a}, \mathbf{b}] \in \mathbb{R}^{m \times 2}$ .
$\mathcal{X}$	Finite vocabulary of words $x$ .
$\mathcal{Y}(x)$	Finite set of trees $y$ that are compatible with $x$ .
$\mathcal{V}(x)$	Finite set of labeled spans $v$ over $x$ .
$X, Y, \dots$	Random variables with sample spaces $\mathcal{X}, \mathcal{Y}, \dots$
$x$	A word from $\mathcal{X}$ , outcome of random variable $X$ .
$y$	A tree from $\mathcal{Y}(x)$ , outcome of random variable $Y$ .
$x_1^m$	A sequence of words $\langle x_1, \dots, x_m \rangle$ from $\mathcal{X}^m$ .
$x_{<i}$	The sequence $x_1^{i-1}$ preceding $x_i$ .
$P_X$	Probability distribution.
$p_X$	Probability mass function.
$p(x)$	Probability $P(X = x)$ .
$p_\theta, q_\lambda$	Probability mass functions emphasizing parameters.
$\mathbb{E}[g(X)]$	Expectation of $g(X)$ with respect to distribution $P_X$ .
$H(P_X)$	Entropy of random variable $X$ with distribution $P$
$\Lambda$	Finite set of nonterminal labels.
$A, B, \dots$	Nonterminal labels from $\Lambda$ .
$S^\dagger$	Special root label not in $\Lambda$ .

# 1

---

## Introduction

---

This thesis investigates the question: *What are effective ways of incorporating syntactic structure into a neural language model?*

We study a class of neural language models that explicitly model the hierarchical syntactic structure in addition to the sequence of words (Dyer *et al.*, 2016; Buys and Blunsom, 2015b; Buys and Blunsom, 2018). These models merges generative transition-based parsing with recurrent neural networks in order to model sentences together with their latent syntactic structure. The syntactic structure that decorates the words can be latent, and marginalized over, or can be given explicitly, for example as the prediction of an external parser. Although these are fundamentally joint model, they can be evaluated as regular language models (modeling only words) by (approximate) marginalization of the syntactic structure. In the case of the RNNG (Dyer *et al.*, 2016), exact marginalization is intractable due to the parametrization of the statistical model, but importance sampling provides an effective approximate method. An externally trained discriminative parser is used to obtain proposal samples. Other models provide exact marginalization, but this typically comes at the cost of a less expressive parametrization, for example one in which the features cannot be structure-dependent (Buys and



Blunsom, 2018).

In this thesis I study the RNNG (Dyer *et al.*, 2016) and investigate:

**The approximate marginalization** I propose an alternative proposal distribution and investigate the impact.

- I propose a new discriminative chart-based neural parser that is trained with a global, Conditional Random Field (CRF), objective. The parser is an adaptation of the minimal neural parser proposed in Stern *et al.* (2017a), which is trained with a margin-based objective.
- This contrast with the choice of Dyer *et al.* (2016) for a transition-based parser as proposal, a discriminatively trained RNNG.
- We posit that a globally trained model is a better proposal distribution than a locally trained transition based model: a global model has ready access to competing analyses that can be structurally dissimilar but close in probability, whereas we hypothesize that a locally trained model is prone to produce locally corrupted structures that are nearby in transition-space. In a transition based parser more diverse samples can be obtained by flattening the transition distributions, which causes the model to be less confident in its predictions. The downside is that the model now explores parts of the probability space which it has not encountered during training.

**Semi-supervised training by including unlabeled data** To make joint models competitive language models they need to make use of the vast amounts of unlabeled data that exists.

- A major drawback of these syntactic language models is that they require annotated data to be trained, and precious little of such data exists.
- We extend the training to the unsupervised domain by optimizing a variational lower bound on the marginal probabilities that jointly

optimizes the parameters of proposal model ('posterior' in this framework) with the joint model.

- We obtain gradients for this objective using the score function estimator (Fu, 2006), also known as REINFORCE (Williams, 1992), which is widely used in the field of deep reinforcement learning, and we introduce an effective baseline based on argmax decoding (Rennie *et al.*, 2017), which significantly reduces the variance in this optimization procedure.
- Our CRF parser particularly excels in the role of posterior thanks to the independence assumptions that allow for efficient exact computation of key quantities: the entropy term in the lower bound can be computed exactly using Inside-Outside algorithm, removing one source of variance from the gradient estimation, and the argmax decoding can be performed exactly thanks to Viterbi, making the argmax baseline even more effective.

**Alternative, simpler, models** There are alternatives to the methods that this thesis investigates.

- Multitask learning of a neural language model with a syntactic side objective is a competitive and robust alternative method to infuse neural language models with syntactic knowledge.
- Training the syntactic model on data that mixes gold trees with predicted 'silver' trees for unlabeled data is a competitive and robust alternative to fully principled semi-supervised learning.
- We propose a simple multitask neural language model that predicts labeled spans from the RNN hidden states, using a feature function identical to that used in the CRF parser.
- We consider these alternatives in order to quantify significance of the latent structure and the semisupervised training as measured by some external performance metric.

**Targeted syntactic evaluation** TBA

# 2

---

## Background

---

This section provides background on the four topics that are combined in this thesis: syntax, parsing, language modelling, and neural networks.

### 2.1 Syntax

This section introduces some concepts from syntax that are relevant for this thesis. In particular, we introduce the notion of a *constituent*, and the hierarchical organization of constituents in a sentence as described by phrase-structure grammars. The aim is to provide a succinct and compelling answer to the question: *Why should we care about constituency structure when modelling language?*

The exposition primarily follows Huddleston and Pullum (2002), a well established reference grammar of the English language that is relatively agnostic with respect to theoretical framework, with some excursions into Carnie (2010) and Everaert *et al.* (2015), which are less language-specific but rooted more in a particular theoretical framework<sup>1</sup>.

We take the following three principles from Huddleston and Pullum (2002) as guiding:

---

<sup>1</sup>Broadly subsumable under the name *generative grammar*.

1. *Sentences consist of parts that may themselves have parts.*
2. *These parts belong to a limited range of types.*
3. *The constituents have specific roles in the larger parts they belong to.*

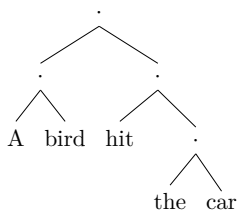
To each principle we now dedicate a separate section.

### 2.1.1 Constituents

Sentences consist of parts that may themselves have parts. The parts are groups of words that function as units and are called *constituents*. Consider the simple sentence *A bird hit the car*. The immediate constituents are *a bird* (the subject) and *hit the car* (the predicate). The phrase *hit the car* can be analyzed further as containing the constituent *the car*. The ultimate constituents of a sentence are the atomic words, and the entire analysis is called the constituent structure of the sentence. This structure can be indicated succinctly with the use of brackets

(1)    [ [ A bird ] [ hit [ the car ] ] ]

or less succinctly as a tree diagram Evidence for the existence of such



constituents can be provided by examples such as the following, which are called constituent tests (Carnie, 2010). Consider inserting the adverb *apparently* into our example sentence, indicating the alleged status of the event described in the sentence. In principle there are six positions available for the placement of *apparently* (including before, and after the sentence). However, only three of these placements are actually permissible:<sup>2</sup>

---

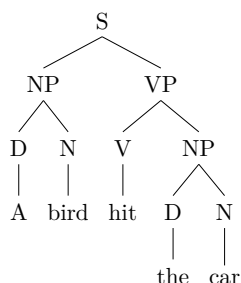
<sup>2</sup>We use an asterisk \* to indicate a sentence that is judged ungrammatical, as is customary in linguistics.

- (2)
- a. *Apparently* a bird hit the car.
  - b. \*An *apparently* bird hit the car.
  - c. A bird *apparently* hit the car.
  - d. \*A bird hit *apparently* the car.
  - e. \*A bird hit the *apparently* car.
  - f. A bird hit the car, *apparently*.

Based on the bracketing that we proposed for this sentence we can formulate a general constraint: the adverb must not interrupt any constituent. Indeed, this would explain why *actually* cannot be placed anywhere inside *hit the car* and not between *a* and *bird*. For full support, typically, results from many more such test are gathered, and in general these tests can be much more controversial than in our simple example (Carnie, 2010).

### 2.1.2 Categories

The constituents of a sentence belong to a limited range of types that form the set of syntactic categories. Two types of categories are distinguished: lexical and phrasal. The lexical categories are also known as part-of-speech tags. A tree can be represented in more detail by adding lexical (D, N, V) and phrasal categories (S, NP, VP). In this



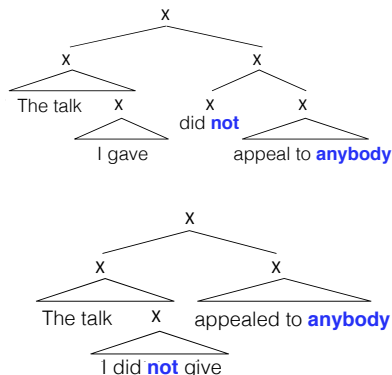
example, the noun (N) *car* is the head of the noun phrase (NP) *the car*, while the head of the larger phrase *hit the car* is the verb (V) *hit*, making this larger constituent a *verb phrase* (VP). The whole combined forms a sentence (S).

### 2.1.3 Hierarchy

The constituents have specific roles in the larger parts they belong to. This structure provides constraints that are not explainable from the linear order of the words themselves (Everaert *et al.*, 2015). Consider an example about the syntactic behaviour of *negative polarity items* (NPIs) from Everaert *et al.* (2015). A negative polarity item is, to first approximation, a word or group of words that is restricted to negative context (Everaert *et al.*, 2015).<sup>3</sup> Take the behaviour of the word *anybody*:

- (3) a. The talk I gave did *not* appeal to *anybody*.  
 b. \*The talk I gave appealed to *anybody*.  
 c. \*The talk I did *not* give appealed to *anybody*.

From sentences (a) and (b) we might formulate the hypothesis that the word *not* must linearly precede the word *anybody*, but a counter example refutes this hypothesis: sentence (c) is also not grammatical. Instead, it



**Figure 2.1:** Hierarchical dependance of negative polarity items. Left shows the word *anybody* in the licensing context of *not*, while right shows the ungrammatical sentence where the word is not. Figure taken without permission from Everaert *et al.*, 2015.

is argued, the constraints that govern this particular pattern depend on

<sup>3</sup>More generally, they are words that need to be licensed by a specific *licencing context* (Giannakidou, 2011).

hierarchical structure: the word *not* must ‘structurally precede’ the word *anybody* (Everaert *et al.*, 2015). Figure shows the constituent structure of both sentences. The explanation goes as follows (put this in my own words): “In sentence 2.1 (a) the hierarchical structure dominating *not* also immediately dominates the hierarchical structure containing *anybody*. In sentence 2.1 (b), by contrast, *not* sequentially precedes *anybody*, but the triangle dominating *not* fails to also dominate the structure containing *anybody*.”

### 2.1.4 Controversy

Theoretical syntax is rife with controversy, and wildly differing viewpoints exist. In fact, for each point made in our short discussion, the exact opposite point has been made as well:

- Work in dependency grammar and other word-based grammar formalisms depart from the idea that lexical relations between individual words are more fundamental than constituents and their hierarchical organization (Tesnière, 1959; Nivre, 2005; Hudson, 2010), and dispenses with the notion of constituents altogether.
- A recurring cause of controversy is the question whether hierarchical structure needs to be invoked in linguistic explanation. That is, whether the kind of analysis in figure ??, with embedded constituents, is really fundamental to language. Frank *et al.* (2012) argue that a shallow analysis of sentences into immediate constituents with linear order but no hierarchical structure<sup>4</sup> is sufficient for syntactic theories, a claim that is vehemently rejected by Everaert *et al.* (2015).
- Research in cognitive neuroscience and psycholinguistics shows that human sentence processing is hierarchical, giving evidence that processing crucially depends on the kind of structures introduced in the sections above (Hale, 2001; Levy, 2008; Brennan *et al.*, 2016). However, research also exists that shows that purely linear

---

<sup>4</sup>A structure reminiscent of that predicted in the task of *chunking*, or shallow parsing.

predictors are sufficient for modeling sentence comprehension, thus showing the exact opposite to be true (Conway and Pisoni, 2008; Gillespie and Pearlmutter, 2011; Christiansen *et al.*, 2012; Gillespie and Pearlmutter, 2013; Frank *et al.*, 2012).

Our work, however, takes a pragmatic position with respect to such questions: syntax, we assume, is whatever our dataset says it is. Whether language is hierarchical or linear is a question that this thesis engages with from a statistical and computational standpoint.

## 2.2 Parsing

Parsing is the task of predicting a tree  $y$  for a sentence  $x$ . Probabilistic parsers solve this problem by learning a probabilistic model  $p$  that describes a probability distribution over *all* the possible parses  $y$  of sentences  $x$ . This distribution quantifies the uncertainty over the syntactic analyses of the sentence, and can be used to make predictions by finding the trees with highest probability. A further benefit of the probabilistic formulation is that the uncertainty over the parses can be determined quantitatively by computing the entropy of the distribution, and qualitatively by obtaining samples from the distribution. This section describes the form that such probabilistic models take, and the data from which they are estimated.

### 2.2.1 Treebank

A treebank is a collection of sentences annotated with their grammatical structure that allows the estimation of statistical parsings models. The Penn Treebank (Marcus *et al.*, 1993) is such a collection, consisting of sentences annotated with their constituency structure. Figure 2.2a shows an example tree from this dataset and figure 2.2b shows the tree after basic processing<sup>5</sup>. Some of the models in this work require trees to be in *normal form*: fully binary, but with unary branches at the terminals. Figure 2.2c show the result of this (invertible) normalization, that is

---

<sup>5</sup>This removes the functional tags and annotation of argument structure introduced in version 2 of the Penn Treebank (Marcus *et al.*, 1994).



obtained by introduction of an empty dummy label  $\emptyset$ . Annotating the labels with left and right endpoints of the words they span results in the tree in figure 2.2d.

A tree can be factorized into its parts: we can think of tree as a set of *labeled spans*, or as a set of *anchored rules*. A labeled span is a triple  $(A, i, j)$  of a syntactic label  $A$  from a labelset  $\Lambda$  together the left and right endpoints  $i, j$  that the label spans. An *anchored rule* is a triple  $(r, i, j)$  or four-tuple  $(r, i, k, j)$ , containing a rule  $r$  in normal form with span endpoints  $i, j$ , and a split-point  $k$  of the left and right child  $r$  is not a lexical rule. For the difference, consider the two representations of the tree in figure 2.2d given in table 2.1. These different factorizations will become relevant in chapter 4.

Labeled spans	Anchored rules
$(S, 0, 10)$	$(S \rightarrow \text{SBAR } \emptyset, 0, 3, 10)$
$(\text{SBAR}, 0, 3)$	$(\text{SBAR} \rightarrow \text{WHNP } S+\text{VP}, 0, 1, 3)$
$(\text{WHNP}, 0, 1)$	$(\text{WHNP} \rightarrow \textit{What}, 0, 1)$
$\vdots$	$\vdots$
$(\emptyset, 9, 10)$	$(\emptyset \rightarrow ., 9, 10)$

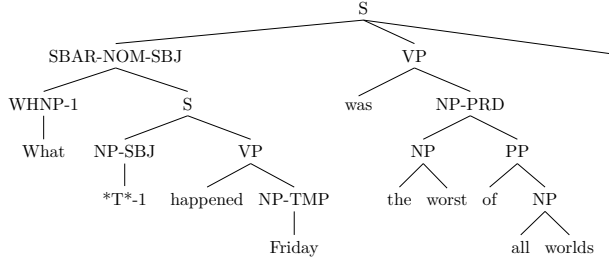
**Table 2.1:** Two conceptions of the tree in 2.2d.

### 2.2.2 Models

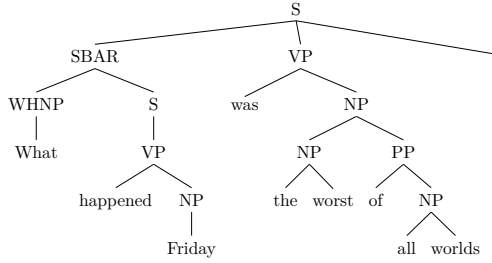
The probabilistic model of a parser can be *discriminative*, describing the conditional probability distribution  $p(y \mid x)$ , or *generative*, describing the joint distribution  $p(x, y)$ . The form that the model takes is largely dictated by the algorithm used to build the parses.

Transition-based methods formulate parsing as a sequence of shift-reduce decisions made by a push-down automaton that incrementally builds a tree. The design of the transition system determines whether the tree is constructed bottom-up, top-down, or otherwise<sup>6</sup>, and determines whether the trees need to be in normal form. The probabilistic model is

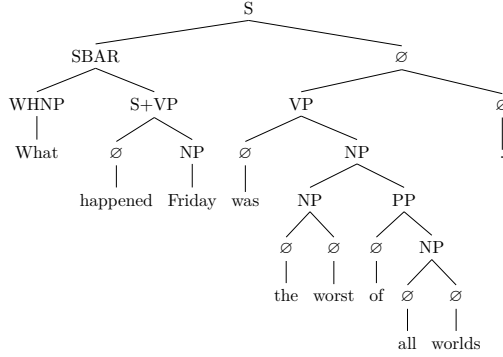
<sup>6</sup>Post-order and pre-order traversal, respectively.



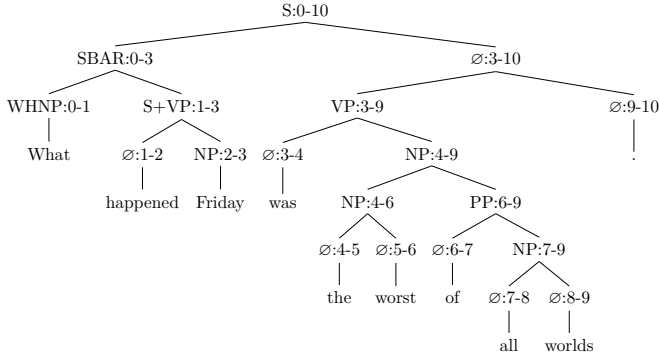
(a) Original Penn Treebank tree.



(b) Function tags and traces removed.



(c) Converted to normal form.



(d) In normal form with spans.

**Figure 2.2:** Converting a treebank tree (withouth part-of-speech tags).

defined over the sequences of actions, and the probability distribution typically factorizes as the product of conditional distributions over the next action given the previous actions. This makes it a directed model that is *locally normalized*.<sup>7</sup>

**Definition 2.1.** A probalistic model  $p(x)$  over sequences  $x \in \mathcal{X}^n$  is *locally normalized* if

$$\begin{aligned} p(x) &= \prod_{i=1}^n p(x_i \mid x_{<i}) \\ &= \prod_{i=1}^n \frac{\Psi(x_{<i}, x_i)}{Z(x_{<i})}, \end{aligned}$$

where  $Z(x_{<i}) = \sum_{x \in \mathcal{X}} \Psi(x_{<i}, x)$  is a local normalizer and  $\Psi$  is a non-negative scoring function:  $\Psi(x_{<i}, x_i) \geq 0$  for all  $x_{<i}$  and  $x$ .

Such a model is discriminative when the actions only build the tree, but can be made generative when word prediction is included in the action set.

Chart based parsing, on the other hand, factorizes trees along their parts, and defines the probabilistic model over the shared substructures. Representing the model as the product of local parts greatly reduces the number of variables required to specify the full distribution, and allows the use of dynamic programming for efficient inference. Discriminative models based on Conditional Random Fields (CRF) (Lafferty *et al.*, 2001) follow this factorization: local functions independently predict scores for the parts, and the product of this score is normalized *globally* by a normalizer that sums over the exponential number structures composable from the parts.

**Definition 2.2.** A probalistic model  $p(x)$  over sequences  $x \in \mathcal{X}^n$  is *globally normalized* if

$$p(x) = \frac{\Psi(x)}{Z},$$

---

<sup>7</sup>With the exception of those approaches that instead define a Conditional Random Field over the action sequences (Andor *et al.*, 2016), in which case the model is defined over the globally normalized action sequences.

where  $Z = \sum_{x \in \mathcal{X}^n} \Psi(x)$  is a *global* normalization term, and  $\Psi(x) \geq 0$  for all  $x$ . To allow efficient computation of the normalizer  $Z$ , the function  $\Psi$  typically factors over parts of  $x$  as  $\Psi(x) = \prod_{a=1}^A \psi(x_a)$ .

Generative chart-based models, instead, estimate rule probabilities directly from a treebank. The probability of a tree is computed directly as the product of the probabilities of the rules that generate it, thus requiring no normalization.<sup>8</sup>

Either method have their advantages and disadvantages. The transition-based methods are fast, running in time linear in the sequence length, and allow arbitrary features that can condition on the entire sentence and the partially constructed derivation. However, directed models with local normalization are known to suffer from the problem of label-bias (Lafferty *et al.*, 2001), and conditioning on large parts of the partial derivations prohibits exact inference, so that decoding can be done only approximately, either with greedy decoding or with beam search. The chart based methods, on the other hand, allow efficient exact inference, and the global training of the discriminative models [...something about learning from all substructures...]. The method is much slower, however, running in time cubic in the sentence length for normal form trees<sup>9</sup> and linear in the size of the grammar<sup>10</sup>, and the strong factorization of the structure, which makes the exact inference tractable, also means that features can condition only on small parts instead of large substructures.

A general challenge for greedy transition is that their training is highly local: they are only exposed to the individual paths that generate the example tree, which is a mere fraction of all the possible paths that the model is defined over.<sup>11</sup> Compare this with a globally normalized

---

<sup>8</sup>In fact, this makes generative chart-based models directed, locally normalized models: they generate trees top-down by expanding localized normally rules.

<sup>9</sup>And slower for trees that are not in normal form.

<sup>10</sup>Giving a total time complexity of  $O(n^3|G|)$ .

<sup>11</sup>One way to answer to this challenge is to use a dynamic oracle during training (Goldberg and Nivre, 2013), also called exploration (Ballesteros *et al.*, 2016; Stern *et al.*, 2017a), which allow the model to explore paths that deviate from the single gold path in a principled way. The approach can be considered an instance of imitation learning (Vlachos, 2013; He *et al.*, 2012). In constituency parsing, dynamic oracles can produce substantial improvements performance (Ballesteros *et al.*, 2016), but they must be custom designed for each transition system (Fried and Klein, 2018).

model, where the factorization over parts lets the model learn from an example tree about all the trees that share substructure.

In this thesis we investigate models where the scoring function  $\Psi$ , or its factorized version  $\psi$ , is implemented with neural networks. Chapter 3 describes a locally normalized transition-based model that has a discriminative and generative formulation, and chapter 4 introduces a globally normalized, discriminative chart-based parser.

### 2.2.3 Metrics

The standard evaluation for parsing is the Parseval metric (Black *et al.*, 1991), which measures the number of correctly predicted labeled spans. Let  $\mathcal{R}$  be the set of labeled spans of the gold reference trees, and let  $\mathcal{P}$ . The metric is defined as the harmonic mean, or  $F_1$  measure, of the labelling recall  $R$  and precision  $P$ :

$$F_1 = \frac{2PR}{P + R}. \quad (2.1)$$

The recall is the fraction of reference spans that are predicted correctly

$$R = \frac{|\mathcal{R} \cap \mathcal{S}|}{|\mathcal{R}|},$$

and the precision is the fraction of the source spans that is correct

$$P = \frac{|\mathcal{R} \cap \mathcal{S}|}{|\mathcal{S}|}.$$

The canonical implementation of the Parseval metric is EVALB (Sekine and Collins, 1997).

## 2.3 Language models

A language model is a probability distribution over sequences of words. There are many ways to design such a model, and many datasets to estimate them from. This section focusses on those that are relevant to the models in this thesis.

---

However, we do not consider this directions in this thesis.

### 2.3.1 Models

A language model is a probabilistic model  $p$  that assigns probabilities to sentences, or more formally, to sequences  $x \in \mathcal{X}^*$  of any length over a finite vocabulary. Factorizing the probability of a single sequence  $x \in \mathcal{X}^m$  over its timesteps gives the directed model:

$$p(x) = \prod_{i=1}^m p(x_i \mid x_{<i}). \quad (2.2)$$

This distribution can be approximated by lower order conditionals that condition on smaller, fixed-size, history, by making the Markov assumption assumption that

$$p(x_i \mid x_{<i}) = p(x_i \mid x_{i-j-1}^{i-1}).$$

This is the approach taken by  $n$ -gram language models. The lower order conditionals can be estimated directly by smoothing occurrence counts (Chen and Goodman, 1999; Kneser and Ney, 1995), or they can be estimated by locally normalized scores  $\Psi(x_{i-j-1}^{i-1}, x_i)$ , given by a parametrized function  $\Psi$ . This function can be log-linear or a non-linear neural network (Rosenfeld, 1996; Bengio *et al.*, 2003). The lower order approximation can also be dispensed with, making the model more expressive, but the estimation problem much harder. Language models based on recurrent neural networks (RNNs) follow this approach by using functions that compute scores  $\Psi(x_{<i}, x_i)$  based on the entire history (Mikolov *et al.*, 2010). These models, and in particular the approaches based on the LSTM variant of the RNN, have shown to be remarkably effective and substantially improve over the above methods to give the current state of the art (Zaremba *et al.*, 2014; Jozefowicz *et al.*, 2016).<sup>12</sup> Other neural network architectures based on convolutions

---

<sup>12</sup>Although recent work shows that the effective memory depth of RNN models is much smaller than the unbounded history suggests: Chelba *et al.* (2017) show that, in the perplexity they assign to test data, an RNN with unbounded history can be approximated very well by an RNN with bounded history. To be precise: a neural  $n$ -gram language model, with the fixed-size history encoded by a bidirectional RNN, is equivalent to an RNN with unbounded history for  $n = 13$ , and to an RNN that is reset at the start of sentences for  $n = 9$ . This finding is consistent accross dataset sizes.

(Kalchbrenner *et al.*, 2014) and stacked feedforward networks with ‘self-attention’ (Vaswani *et al.*, 2017) have been successfully applied to language modelling with unbounded histories as well.

Alternatively, a language model can be obtained by marginalizing a structured latent variable  $h$  in a joint model  $p(x, h)$ :

$$p(x) = \sum_{h \in \mathcal{H}} p(x, h). \quad (2.3)$$

The structure of  $h$  allows this joint distribution to be factorized in a ways very much unlike that of equation 2.2. Such language models are defined for example by a Probabilistic Context Free Grammar (PCFG), in which case  $h$  is a tree, and a Hidden Markov Model (HMM), in which case  $h$  is a sequence of tags. The strong independence assumptions of these models allows the marginalization to be computed efficiently, but also disallows the models to capture dependencies in  $x$ , which is precisely what we want from a language model. The recurrent neural network grammar (RNNG) model introduced in chapter 3 also defines a language model through a joint distribution, but that model is factorized as a sequential generation process over both the latent structure  $h$  and the observed sequence  $x$ , which makes it a competitive language model, but at the price of losing efficient exact marginalization.

Some other language models that incorporate syntax: language models obtained from top-down parsing with a PCFG (Roark, 2001); syntactic extensions of  $n$ -gram models with count-based and neural network estimation (Chelba and Jelinek, 2000; Emami and Jelinek, 2005); and an method that is reminiscent of  $n$ -gram methods, but based on arbitrary overlapping tree fragments (Pauls and Klein, 2012).

### 2.3.2 Data

Language models enjoy the benefit that they require no labeled data; any collection of tokenized text can be used for training. In this work we focus on English language datasets. The tokens from Penn Treebank have long been a popular dataset for this task. More recently has seen the introduction of datasets of much greater size, such as the One Billion Word Benchmark (Chelba *et al.*, 2013) that consists of news articles, and datasets that focus on long-distance dependencies, such as

the Wikitext datasets (Merity *et al.*, 2016) that consists of Wikipedia articles grouped into paragraphs.

### 2.3.3 Metrics

The standard metric to evaluate a language model is the *perplexity* that it assigns to held out data. The lower the perplexity, the better the model. Perplexity is an information theoretic metric that corresponds to an exponentiated estimate of the model’s entropy, measured in nats, and was introduced for this purpose by Jelinek (1997)<sup>13</sup>. The metric can be interpreted as the average number of guesses needed by the model to predict each word from its left context.

**Definition 2.3.** The perplexity of a language model  $p$  on a sentence  $x$  of length  $m$  is defined as

$$\exp \left\{ - \frac{1}{m} \log p(x) \right\}.$$

The perplexity on a set of sentences  $\{x^{(1)}, \dots, x^{(N)}\}$  is defined as the exponentiated mean over all words together:

$$\exp \left\{ - \frac{1}{M} \sum_{i=1}^N \log p(x^{(i)}) \right\},$$

where  $M = \sum_{i=1}^N m_i$  is the sum of all the sentence lengths  $m_i$ .

The appeal of perplexity is that it is an aggregate metric, conflating different causes of success when predicting the next word. This conflation is also its main shortcoming, making it hard to determine whether the model has robustly learned high-level linguistic patterns such as those described in syntactic and semantic theories. For this reason, alternative methods have been proposed to evaluate language models: evaluation with adversarial examples (Smith, 2012); prediction of long distance subject-verb agreement (Linzen *et al.*, 2016); and eliciting syntactic acceptability judgments (Marvin and Linzen, 2018). Chapter 6 discusses these alternatives in greater detail, and demonstrates evaluation with the method proposed in (Marvin and Linzen, 2018).

---

<sup>13</sup>According to (Chelba *et al.*, 2017).



## 2.4 Neural networks

In this thesis we use neural networks to parametrize probability distributions. We consider a neural network as an abstraction that denotes a certain type of parametrized differentiable function, and use various kinds. Let  $x$  be a word from a finite vocabulary  $\mathcal{X}$ , and let  $\mathbf{x}$  and  $\mathbf{y}$  be vectors in respectively  $\mathbb{R}^n$  and  $\mathbb{R}^m$ .

**Definition 2.4.** A *word embedding* is vector representation of a word, assigned by an embedding function  $E$  that takes elements from  $\mathcal{X}$  to  $\mathbb{R}^n$ :

$$\mathbf{x} = E(x).$$

The function can be a simple lookup table, or a more elaborate function that depends for example on the orthography of the word.

**Definition 2.5.** A *feedforward neural network* is a parametrized function FFN from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ :

$$\mathbf{y} = \text{FFN}(\mathbf{x}).$$

Internally, the function computes an affine transformation followed by an elementwise application of a nonlinear function, repeatedly. The number of repetitions is referred to as the number of layers of the network.

**Definition 2.6.** A *recurrent neural network* (RNN) is a parametrized function  $\text{RNN}$  that takes a sequence of vectors  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k]$ , each in  $\mathbb{R}^n$ ,<sup>14</sup> and produces a sequence of output vectors  $[\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n]$  each in  $\mathbb{R}^m$ :

$$[\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k] = \text{RNN}(\mathbf{X}).$$

Each vector  $\mathbf{y}_i$  is a function of the vectors  $[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_i]$ , for which reason we refer to the vector  $\mathbf{y}_i$  as a context-dependent *encoding* of the vector  $\mathbf{x}_i$ . An RNN can be applied to the input sequence in reverse. This makes each  $\mathbf{y}_i$  a function of the vectors  $[\mathbf{x}_i, \mathbf{x}_{i+1}, \dots, \mathbf{x}_k]$ . We denote the output of the forward direction with  $\mathbf{f}$  and the output of the backward direction with  $\mathbf{b}$ .

---

<sup>14</sup>Compactly represented as a matrix  $\mathbf{X} \in \mathbb{R}^{n \times k}$ .

**Definition 2.7.** An RNN is *bidirectional* when it combines the output of an RNN that runs in the forward direction, with the output of an RNN that runs in the backward direction. Combining their output vectors by concatenation gives for each position a vector  $\mathbf{h}_i = [\mathbf{f}_i; \mathbf{b}_i]$  that is a function of the entire input sequence.

**Definition 2.8.** An LSTM (Hochreiter and Schmidhuber, [1997](#)) is a particular way to implement the internals of the RNN function. It is the only type of RNN used in this work, and we use the two names exchangeably.

# 3

---

## Recurrent Neural Network Grammars

---

In this chapter we introduce the Recurrent Neural Network Grammar (RNNG), a probabilistic model of sentences with explicit phrase structure, introduced by Dyer *et al.* (2016).

### 3.1 Model

The discriminative RNNG is a discriminative transition based parser that uses regular RNNs to summarize the actions in the history and the words on the buffer into vectors, and uses a special RNN with a syntax-dependent recurrence to obtain a vector representation of items on the the stack. The generative RNNG can be seen as a generative formulation of the RNNG, jointly modelling the words instead of merely conditioning on them, or it can be seen as a structured model of language that predicts words together with their structure. From a parsing perspective, the generative RNNG simply dispends with the discriminative RNNG's buffer to instead predict the words of the tree. As a model of sentences, it can better be understood as kind of structured RNN: it predicts words incrementally, but compresses and labels them recursively whenever they form a complete constituents.

### 3.1.1 Transition system

The RNNG uses a transition system crafted for top-down parsing. The discriminative transition system has three actions:

- OPEN( $X$ ) opens a new nonterminal symbol  $X$  in the partially constructed tree.
- SHIFT moves the topmost word  $x$  from the buffer onto the stack.
- REDUCE closes the open constituent by composing the symbols in it into a single item representing the content of the subtree. This allows nodes with an arbitrary number of children.

How these actions work to build a tree is best illustrated with an example derivation.

**Example 3.1.** (Discriminative transition system) Producing the gold tree of the input sentence *The hungry cat sleeps..* Individual items are

	Stack	Buffer	Action
0		<i>The</i>   <i>hungry</i>   <i>cat</i>   <i>meows</i>   .	OPEN(S)
1	(S	<i>The</i>   <i>hungry</i>   <i>cat</i>   <i>meows</i>   .	OPEN(NP)
2	(S   (NP	<i>The</i>   <i>hungry</i>   <i>cat</i>   <i>meows</i>   .	SHIFT
3	(S   (NP   <i>The</i>	<i>hungry</i>   <i>cat</i>   <i>meows</i>   .	SHIFT
4	(S   (NP   <i>The</i>   <i>hungry</i>	<i>cat</i>   <i>meows</i>   .	SHIFT
5	(S   (NP   <i>The</i>   <i>hungry</i>   <i>cat</i>	<i>meows</i>   .	REDUCE
6	(S   (NP <i>The hungry cat</i> )	<i>meows</i>   .	OPEN(VP)
7	(S   (NP <i>The hungry cat</i> )   (VP	<i>meows</i>   .	SHIFT
8	(S   (NP <i>The hungry cat</i> )   (VP   <i>meows</i>	.	REDUCE
9	(S   (NP <i>The hungry cat</i> )   (VP <i>meows</i> )	.	SHIFT
10	(S   (NP <i>The hungry cat</i> )   (VP <i>meows</i> )   .	.	REDUCE
11	(S (NP <i>The hungry cat</i> ) (VP <i>meows</i> ) .)		

separated by the midline symbol.

The actions are constrained in order to only derive well-formed trees:

- SHIFT can only be executed if there is at least one open nonterminal (all words must be under some nonterminal symbol).
- OPEN( $X$ ) requires there to be words left on the buffer (all constituents must eventually hold terminals).
- REDUCE requires there to be at least one terminal following the open nonterminal, and the topmost nonterminal can only be closed when the buffer is empty (all nodes must end under a single root node).

The generative transition system is derived from this system by replacing the SHIFT action, which moves the word  $x$  from the buffer into the open constituent on the stack, with an action that *predicts* the word:

- $\text{GEN}(x)$  predicts that  $x$  is the next word in the currently open constituent, and puts this word on the top of the stack.

The buffer is dispensed with and is replaced by a similar structure that records the sequence words predicted so far.

**Example 3.2.** (Generative transition system) Generating the sentence of example 3.1 with its gold tree.

	Stack	Terminals	Action
0			OPEN(S)
1	(S		OPEN(NP)
2	(S   (NP		GEN( <i>The</i> )
3	(S   (NP   <i>The</i>	<i>The</i>	GEN( <i>hungry</i> )
4	(S   (NP   <i>The</i>   <i>hungry</i>	<i>The</i>   <i>hungry</i>	GEN( <i>cat</i> )
5	(S   (NP   <i>The</i>   <i>hungry</i>   <i>cat</i>	<i>The</i>   <i>hungry</i>   <i>cat</i>	REDUCE
6	(S   (NP <i>The hungry cat</i> )	<i>The</i>   <i>hungry</i>   <i>cat</i>	OPEN(VP)
7	(S   (NP <i>The hungry cat</i> )   (VP	<i>The</i>   <i>hungry</i>   <i>cat</i>	GEN( <i>meows</i> )
8	(S   (NP <i>The hungry cat</i> )   (VP   <i>meows</i>	<i>The</i>   <i>hungry</i>   <i>cat</i>   <i>meows</i>	REDUCE
9	(S   (NP <i>The hungry cat</i> )   (VP <i>meows</i> )	<i>The</i>   <i>hungry</i>   <i>cat</i>   <i>meows</i>	GEN(.)
10	(S   (NP <i>The hungry cat</i> )   (VP <i>meows</i> )   .	<i>The</i>   <i>hungry</i>   <i>cat</i>   <i>meows</i>   .	REDUCE
11	(S (NP <i>The hungry cat</i> ) (VP <i>meows</i> ) .)	<i>The</i>   <i>hungry</i>   <i>cat</i>   <i>meows</i>   .	

### 3.1.2 Model

Fundamentally, the model is a probability distribution over transition actions  $a = (a_1, \dots, a_T)$  that generates a tree  $y$ . Conditionally, given a sequence of words  $x$ , in the discriminative model, and jointly predicting  $x$  in the generative model. Put simply, the model is thus defined as

$$p(a) = \prod_{t=1}^T p(a_t \mid a_{<t}), \quad (3.1)$$

where in the discriminative case  $p(a) = p(y \mid x)$  and in the generative model  $p(a) = p(x, y)$ .

The exact model however is slightly more complicated, a consequence of the difference between the discriminative and the generative actions, and a consequence of practical concerns regarding the implementation.

To define the model precisely we need to introduce some things. First we define the set of discriminative actions as

$$A_{\mathcal{D}} = \{\text{SHIFT}, \text{OPEN}, \text{REDUCE}\}, \quad (3.2)$$

and the set of generative actions as

$$A_{\mathcal{G}} = \{\text{GEN}, \text{OPEN}, \text{REDUCE}\}. \quad (3.3)$$

We define a finite set of nonterminal symbols  $N$  and a finite alphabet of words  $\Sigma$ . Then in the discriminative model  $a$  is an element of  $A_{\mathcal{D}}^T$ , and in the generative model  $a$  is an element of  $A_{\mathcal{G}}^T$ , both with the restriction that the actions form a valid tree  $y$ . The sequence of nonterminals  $\mathbf{n}$  in  $N^K$  is the sequence of nonterminal nodes obtained from  $y$  by pre-order traversal. An a sentence  $x$ , finally, is an element of  $\Sigma^N$ .

Finally, let  $\mathbf{1}_{\{a_t = \text{OPEN}\}}$  be the indicator function for the event that the action  $a_t$  is to open a new nonterminal, and similarly let  $\mathbf{1}_{\{a_t = \text{GEN}\}}$  be the indicator function for the event that the action  $a_t$  is to generate a word. Then we can introduce two functions that map between sets of indices to indicate the number of times a particular action has been taken at each time step:

$$\mu : [T] \rightarrow [M] : \mu(t) = \sum_{i=0}^{t-1} \mathbf{1}_{\{a_i = \text{OPEN}\}},$$

and

$$\nu : [T] \rightarrow [N] : \nu(t) = \sum_{i=0}^{t-1} \mathbf{1}_{\{a_i = \text{GEN}\}}.$$

We are now in the position to write down the exact models. Let  $a$  be a sequence from  $A_{\mathcal{D}}^T$ . Then the model for the discriminative RNNG is

$$p(y \mid x) = p(a \mid x) = \prod_{t=1}^T p(a_t \mid x, a_{<t}), \quad (3.4)$$

where

$$p(a_t \mid x, a_{<t}) = \begin{cases} p(n_{a_t} \mid x, a_{<t}) p(n_{\mu(t)} \mid x, a_{<t}) & \text{if } a_t = \text{OPEN} \\ p(n_{a_t} \mid x, a_{<t}) & \text{otherwise.} \end{cases} \quad (3.5)$$

Let  $a$  be a sequence from  $A_G^T$ , which include the actions that generate words<sup>1</sup>, then the model for the generative RNNG is

$$p(y \mid x) = p(a) = \prod_{t=1}^T p(a_t \mid a_{<t}), \quad (3.6)$$

where

$$p(a_t \mid a_{<t}) = \begin{cases} p(a_t \mid a_{<t})p(n_{\mu(t)} \mid a_{<t}) & \text{if } a_t = \text{OPEN} \\ p(a_t \mid a_{<t})p(x_{\nu(t)} \mid a_{<t}) & \text{if } a_t = \text{GEN} \\ p(a_t \mid a_{<t}) & \text{otherwise.} \end{cases} \quad (3.7)$$

The probabilities over next actions at time step  $t$  are given by classifiers on a vector  $\mathbf{u}_t$ <sup>2</sup>, which represents the parser configuration at that timestep:

$$p(a_t \mid a_{<t}) \propto \exp \text{FFN}_{\alpha}(\mathbf{u}_t) \quad (3.8)$$

$$p(n_{\mu(t)} \mid a_{<t}) \propto \exp \text{FFN}_{\beta}(\mathbf{u}_t) \quad (3.9)$$

$$p(x_{\nu(t)} \mid a_{<t}) \propto \exp \text{FFN}_{\gamma}(\mathbf{u}_t) \quad (3.10)$$

where  $\{\alpha, \beta, \gamma\}$  are separate sets of parameters. The computation of the vector  $\mathbf{u}_t$  is described in the next section.

**Remark 3.1.** We could have defined the actions as

$$A_{\mathcal{D}} = \{\text{REDUCE}, \text{SHIFT}\} \cup \{\text{OPEN}(n) \mid n \in N\},$$

and

$$A_G = \{\text{REDUCE}\} \cup \{\text{OPEN}(n) \mid n \in N\} \cup \{\text{GEN}(x) \mid x \in \Sigma\},$$

and defined  $p(a)$  as in 3.1. However, in the case of the generative model this is particularly inefficient from a computational perspective. Note that the set  $\Sigma$  is generally very very large<sup>3</sup>, and the normalization in 3.8 requires a sum over all actions while a large number of the actions do not

<sup>1</sup>Note that under this action set  $p(a_t \mid a_{<t}, x_{<t}) = p(a_t \mid a_{<t})$ , given the fact that the words in  $x_{<t}$  are contained in  $a_{<t}$ .

<sup>2</sup>For brevity we omit the conditioning on  $x$ , which was redundant already in the case of the generative model.

<sup>3</sup>On the order of tens or hundreds of thousands.

generate words. Besides, the presentation in 3.6 and 3.7 is conceptually cleaner: first choose an action, then, if required, choose the details of that action. For these reasons we opt for the two-step prediction. For consistency we extend this modelling choice to the discriminative RNNG. And although it appears that Dyer *et al.* (2016) model the sequences according to 3.1, followup work takes our approach and models the actions of the generative RNNG as in equation 3.7 (Hale *et al.*, 2018).

### 3.2 Parametrization

The transition probabilities are computed from the vector  $\mathbf{u}_t$  that summarizes the parser’s entire configuration history at time  $t$ . This vector is computed incrementally and in a syntax-dependent way. It is defined as the concatenation of three vectors, each summarizing one of the three datastructures separately:

$$\mathbf{u}_t = [\mathbf{s}_t, \mathbf{b}_t, \mathbf{h}_t],$$

Here,  $\mathbf{s}_t$  represents the stack,  $\mathbf{b}_t$  represents the buffer, and  $\mathbf{h}_t$  represents the history of actions.

The vectors  $\mathbf{b}_t$  and  $\mathbf{h}_t$  are computed each with a regular RNN: the history vector is computed in the forward direction, and the buffer is encoded in the backward direction to provide a lookahead. The vector  $\mathbf{s}_t$  represents the partially constructed tree that is on the stack, and its computation depends on this partial structure by using a structured RNN that encodes the tree in top-down order while recursively compressing constituents whenever they are completed.

#### 3.2.1 Stack encoder

The stack RNN computes a representation based on incoming nonterminal and terminal symbols while these are respectively opened and shifted, as a regular RNN would, but rewrites this history whenever the constituent that they form is closed. Whenever a REDUCE action is predicted, the RNN rewinds its hidden state to before the constituent was opened; the items making up the constituent are composed into a single vector by a composition function; and this composed vector



is then fed into the rewound RNN as if the subtree where a single input. The closed constituent is now a single item represented by a single vector. Due to the nested nature of constituents this procedure recursively compresses subtrees.

**Example 3.3.** (Composition function) Consider the state of the parser in example 3.1 at step 5, when the stack contains the five items

$$(S \mid (NP \mid The \mid hungry \mid cat).$$

Each first represented by an embedding vector, and then encoded by an RNN in the regular way. A REDUCE action is now predicted: the items are popped from the stack until the first open bracket is popped, which in this example is the item that opens the NP bracket. This process also rewinds the RNN state to before the bracket was opened, which in this example brings the RNN back to its state at step 1. The composition now computes a vector representation for the 5 popped items, returning a single representation for the item

$$(NP \ The \ hungry \ cat).$$

The RNN is fed this input reduced input, resulting in the encoding of stack at step 6, and finalizing the reduction step. This process is recursive: consider the reduction that takes the stack from the five items

$$(S \mid (NP \mid The \mid (ADJP \ very \ hungry) \mid cat$$

to the two items

$$(S \mid (NP \ The \ (ADJP \ very \ hungry) \ cat).$$

The items in the constituent (ADJP *very hungry*) have already been composed into a single item, and now it takes part in the composition at a higher level.

### 3.2.2 Composition function

Two kinds of functions have been proposed to for the composition described above: the original function based on a bidirectional RNN (Dyer *et al.*, 2016), and a more elaborate one that additionally incorporates

an attention mechanism (Kuncoro *et al.*, 2017). Both methods encode the individual items that make up the constituent with a bidirectional RNN but while the simpler version merely concatenates the endpoint vectors, the attention based method computes a convex combination of all these vectors, weighted by predicted attention weights. Kuncoro *et al.* (2017) show that the attention-based composition performs best and so we only consider this function.

### 3.3 Training

The discriminative and generative model are trained to maximize the objective

$$\mathcal{L}(\theta) = \sum_{(x,y) \in \mathcal{D}} p_{\theta}(y \mid x),$$

respectively

$$\mathcal{L}(\theta) = \sum_{(x,y) \in \mathcal{D}} p_{\theta}(x, y),$$

by gradient-based optimization on the parameters  $\theta$ .

### 3.4 Inference

The two formulations of the RNNG have complementary applications: both models can be used for parsing, but the generative model can additionally be used as a language model. How these problems can be solved is the topic of this section.

#### 3.4.1 Discriminative model

Parsing a sentence  $x$  with the discriminative model corresponds to solving the following search problem of finding the maximum *a posteriori* (MAP) tree

$$\hat{y} = \arg \max_{y \in \mathcal{Y}(x)} p_{\theta}(y \mid x).$$

Solving this exactly is intractable due to the parametrization of the model. At each timestep, the model conditions on the entire history

derivation history which excludes the use of dynamic programming to solve this efficiently. Instead we rely on an approximate search strategy. There are two common approaches for this: greedy decoding and beam search. We only focus on the first. Greedy decoding is precisely what it suggests: at each timestep we greedily select the best local decision

$$a_t^* = \arg \max_a p_\theta(a | a_{<t}^*).$$

The predicted parse is then the approximate MAP tree  $y^*$  constructed by the sequence  $a^* = (a_1^*, \dots, a_m^*)$ .

### 3.4.2 Generative model

Given a trained generative model  $p_\theta$  we are interested in solving the following two problems: parsing a sentence  $x$

$$\hat{y} = \arg \max_{y \in \mathcal{Y}(x)} p_\theta(x, y),$$

and computing its marginal probability

$$p(x) = \sum_{y \in \mathcal{Y}(x)} p_\theta(x, y).$$

Either inference problem is intractable as either problem would require exhaustive enumeration of and computation over all possible action sequences. Luckily, both can be effectively approximated with the same method: importance sampling using a conditional proposal distribution  $q_\lambda(y | x)$  (Dyer *et al.*, 2016).

**Approximate marginalization** The proposal distribution allows us to rewrite the marginal probability of a sentence as an expectation under this distribution:

$$\begin{aligned} p(x) &= \sum_{y \in \mathcal{Y}(x)} p_\theta(x, y) \\ &= \sum_{y \in \mathcal{Y}(x)} q_\lambda(y | x) \frac{p_\theta(x, y)}{q_\lambda(y | x)} \\ &= \mathbb{E}_q \left[ \frac{p_\theta(x, y)}{q_\lambda(y | x)} \right] \end{aligned}$$

This expectation can be approximated with a Monte Carlo estimate

$$\mathbb{E}_q \left[ \frac{p_\theta(x, y)}{q_\lambda(y|x)} \right] \approx \frac{1}{K} \sum_{i=1}^K \frac{p_\theta(x, y_i)}{q_\lambda(y_i|x)}, \quad (3.11)$$

using proposal samples  $y_i$  sampled from the proposal model  $q_\lambda$  conditioning on  $x$ .

**Approximate MAP tree** To approximate the MAP tree  $\hat{y}$  we use the same set of proposal samples as above, and choose the tree  $y$  from the proposal samples that has the highest probability under the joint model  $p_\theta(y, x)$ .

**Proposal distribution** In order to avoid division by zero, the proposal distribution must satisfy the property that for all  $x$  and  $y \in \mathcal{Y}(x)$

$$p(x, y) > 0 \Rightarrow q(y|x) > 0.$$

We additionally want that samples can be obtained efficiently, and that their conditional probability can be computed. All these requirements are met by the discriminative RNNG: samples can be obtained by ancestral sampling over the transition sequences. For this reason Dyer *et al.* (2016) use this as their proposal. However, any discriminatively trained parser that meets these requirements can be used alternatively, and in chapter 4 we will propose such an alternative.

## 3.5 Experiments

We perform three types of experiments with the RNNG:

- We reproduce the parsing f-scores and perplexities from (Dyer *et al.*, 2016), and some more.
- We evaluate ‘how good the model is’ as a sampler.

### 3.5.1 Supervised model

We investigate the following.

- We train with standard hyperparameter settings and optimizer, and replicate the original results. We will get a little lower with the discriminative model because we do not use tags.
- We evaluate F-score with 100 samples (as many proposal trees as possible).
- We evaluate perplexity with varying number of samples: 1 (argmax), 10, 20, 50, 100 (default). The perplexity evaluation with the argmax prediction gives an impression of the uncertainty in the model (Buys and Blunsom, 2018).

### 3.5.2 Sampler

We investigate the following:

- We assess the conditional entropy of the model. This is most quantitative. Recall that conditional entropy is defined as

$$H(Y | X) = \sum_{x \in \mathcal{X}} p_X(x) H(Y | X = x), \quad (3.12)$$

where

$$H(Y | X = x) = - \sum_{y \in \mathcal{Y}} p_{Y|X}(y | x) \log p_{Y|X}(y | x). \quad (3.13)$$

We estimate the quantity  $H(Y | X = x)$  with the model samples. We estimate the quantity  $H(Y | X)$  by a sum over the development dataset. For the probabilities  $p_X(x)$  we use the marginalized probabilities of the joint RNNG (with samples from the discriminative parser  $p_{Y|X}$ ).

- We assess for some cherry picked sentences. This is more qualitative. These sentences should be difficult or ambiguous. Or they can be ungrammatical when taken from the syneval dataset. We can evaluate their entropy, and the diversity of samples, for example to see if there are clear modes. We can make violinplots of the probabilities of the samples. We can compute the f-scores of the samples compared with the argmax tree.

## 3.6 Related work

### 3.6.1 Generative parsing

- Generative dependency parsing and language modelling (Titov and Henderson, 2007; Buys and Blunsom, 2015a; Buys and Blunsom, 2015b; Buys and Blunsom, 2018)
- Top-down parsing and language modelling (Roark, 2001).

### 3.6.2 Syntax

RNNGs learn about syntax beyond their supervision. The attention mechanism in the composition function learns a type of ‘soft’ head-rules<sup>4</sup>, in which most of the attention weight is put on the item that linguists consider the syntactic head of such constituents. Another finding shows that when trained on unlabeled trees, the RNNG learns representations for constituents that cluster according to their withheld gold label (Kuncoro *et al.*, 2017). Additionally RNNGs are much better at a long-distance subject-verb agreement task than LSTMs (Linzen *et al.*, 2016; Kuncoro *et al.*, 2018), which advantage they owe to the composition function that by repeatedly compressing intervening structure decreases the distance between the two words at stake.

### 3.6.3 Cognition

RNNGs can tell us something about our brains. Psycholinguistic research has shown that top-down parsing is a cognitively plausible parsing strategy (Brennan *et al.*, 2016), and recently, RNNGs have been shown to be particularly good statistical predictors for human sentence comprehension (Hale *et al.*, 2018). In this experiment, the sequential word-probabilities derived from a generative RNNG<sup>5</sup> provide a per-word complexity metric that predict human reading difficulty well. Much better at least than predictions from the word probabilities obtained from a purely sequential RNN language model.

---

<sup>4</sup>A head is the lexical item in a phrase that determines the syntactic category of that phrase, *cf.* the background.

<sup>5</sup>Obtained by using ‘word-synchronous beam-search’ (Stern *et al.*, 2017b).

# 4

---

## Conditional Random Field parser

---

This chapter introduces a neural Conditional Random Field (CRF) parser that can be used as an alternative proposal distribution for the approximate marginalization of the RNNG. The model is a span-factored CRF that independently predicts scores for labeled spans over the sentence using neural networks. The scores then interact in a tree-structured dynamic program, giving a compact description of the probability distribution over parses. This approach combines the efficient exact inference of chart-based parsing, the rich nonlinear features of neural networks, and the global training of a CRF. The chart-based approach enables efficient exact inference allowing for exact decoding and global sampling while the neural features can be complex and can condition on the entire sentence. The CRF training objective [...something about receiving information about all substructures...]. The parser is an adaptation of the chart-based parser introduced in Stern *et al.* (2017a) to global CRF training. In Stern *et al.* (2017a) the model is trained with a margin-based objective, but given our interest in this model as a proposal distribution we require probabilistic training.

In this chapter:

- I present the parser and describe how it is an adaptation of earlier

log-linear and neural crf parsing (Finkel *et al.*, 2008; Durrett and Klein, 2015) and the maximum margin trained parser of Stern *et al.* (2017a).

- I show how the several inference problems of interest can be solved exactly: prediction, entropy, sampling.
- We train the parser in a supervised fashion and show it's adequacy as a parser.
- We investigate the kind of samples that are obtained from this parser, and evaluate how they affect the approximate inference in the RNNG.

## 4.1 Model

The model is a CRF factored over labeled spans. Let  $x$  be a sentence, and  $y$  a tree from  $\mathcal{Y}(x)$ . We define a function  $\Psi$  that assigns nonnegative scores  $\Psi(x, y)$  and let the probability of a tree be its globally normalized score

$$p(y \mid x) = \frac{\Psi(x, y)}{Z(x)}, \quad (4.1)$$

where

$$Z(x) = \sum_{y \in \mathcal{Y}(x)} \Psi(x, y)$$

is the normalizer, or partition function, that sums over the exponential number of trees available for  $x$ .

To allow efficient computation of the normalizer, we let the scoring function  $\Psi$  factorize over the parts of  $y$ . We consider a tree as a set of labeled spans  $y_a = (A, i, j)$ , thus  $y = \{y_a\}_{a=1}^A$ , where  $(A, i, j)$  indicates that  $y$  contains a label  $A$  spanning the words  $\langle x_{i+1}, \dots, x_j \rangle$ . The value  $\Psi(x, y)$  is then defined as the product of the nonnegative potentials  $\psi(x, y_a)$  as

$$\Psi(x, y) = \prod_{a=1}^A \psi(x, y_a). \quad (4.2)$$



The function  $\psi$ , thus, scores each labeled span  $y_a$  separately, but conditional on the entire sentence.

The above model is your typical constituency parsing CRF (Finkel *et al.*, 2008; Durrett and Klein, 2015), but the factorization over labeled spans, however, was first introduced by Stern *et al.* (2017a). Factorizing a tree over labeled spans makes an even stronger assumption than that typically made of factorizing over anchored rules<sup>1</sup>, such as is done in the earlier work. By factorizing over labeled spans, the potential function  $\psi$  has no access to information about the direct substructure under the node, such as the child nodes and their split point, and the function can thus rely less on the (local) tree structure and must thus rely more on the surface features of the input sentence. This factorization will, however, greatly reduce the size of the state-space of the dynamic programs, speeding up training and inference, and the burden of the feature function will be carried by a rich neural network parametrization. Together this will make the parser fast yet effective. This will be made clear in section 4.3 on inference.

## 4.2 Parametrization

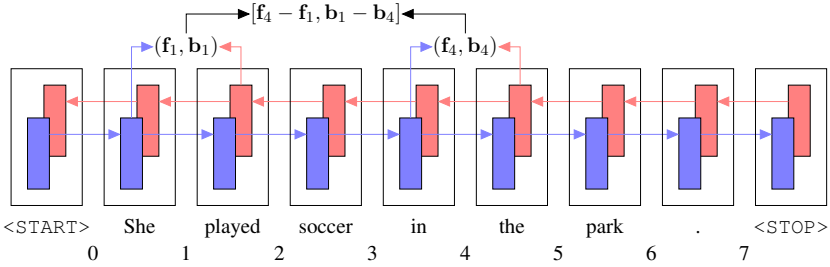
The scoring function  $\psi$  is implemented with neural networks following Stern *et al.* (2017a). Again, let  $y_a$  denote a labeled span  $(A, i, j)$  in a tree  $y$ , and let  $\psi(x, y_a) \geq 0$  be the score of that labeled span for sentence  $x$ . These local potentials can only make minimal use of structural information but they can depend on the entire sentence. This suggests the use of bidirectional RNN encodings. Let  $\mathbf{f}_i$  and  $\mathbf{b}_i$  respectively be the vectors computed by a forward and backward RNN for the word in position  $i$ . The representation of the span  $(i, j)$  is the concatenation of the difference between the vectors on the endpoints of the span:

$$\mathbf{s}_{ij} = [\mathbf{f}_j - \mathbf{f}_i; \mathbf{b}_i - \mathbf{b}_j]. \quad (4.3)$$

The vector  $\mathbf{s}_{ij}$  represents the words  $x_i^j$ , and equation 4.3 is illustrated in figure 4.1. The scores the labels of that position are computed from this vector using a feedforward network with output dimension  $\mathbb{R}^{|\Lambda|}$ ,

---

<sup>1</sup>Compare table 2.1 for the different conceptions of a tree.



**Figure 4.1:** Representation for the span (1, 4) computed from RNN encodings. Taken without permission from Gaddy *et al.* (2018).

and the score of label  $A$  is given by the index corresponding to it:

$$\log \psi(x, y_a) = [\text{FNN}(\mathbf{s}_{ij})]_A, \quad (4.4)$$

where we pretend that  $A$  doubles as an index. This architecture is rightly called minimal, but it works surprisingly well: Stern *et al.* (2017a) experiment with more elaborate functions based on concatenation of vectors (a strict superset of the minus approach) and biaffine scoring (inspired by Dozat and Manning (2016)), but these functions improve marginally, if they do at all.

### 4.3 Inference

Because the model is span-factored it allows efficient inference. In this section we describe efficient solutions to four related problems:

- Compute the normalizer  $Z(x) = \sum_{y \in \mathcal{Y}(x)} \prod_{a=1}^A \psi(x, y_a)$ .
- Find the best parse  $\hat{y} = \arg \max_y p(y | x)$
- Sample a tree  $Y \sim P(Y | X = x)$ .
- Compute the entropy  $H(Y | X = x)$  over parses for  $x$ .

These problems can be solved by instances of the *inside algorithm* and *outside algorithm* (Baker, 1979) with different semirings, an insight we take from semiring parsing (Goodman, 1999). In the following derivations we will make use of the notion of a *weighted hypergraph* as

a compact representation of all parses and their scores (Gallo *et al.*, 1993; Klein and Manning, 2004), and use some of the ideas and notation of *semiring parsing* (Goodman, 1999; Li and Eisner, 2009). First we describe the structure of the parse forest specified by our CRF parser, and then derive the particular form of the inside and outside recursions for this hypergraph from the general formulations. We refer the reader to appendix C for background on these ideas, and the introduction of the notation.

#### 4.3.1 Weighted parse forest

The hypergraph  $G = (\mathcal{V}, \mathcal{E})$  specified by the CRF parser has the following structure.

The set  $\mathcal{V}$  is defined relative to the sentence  $x$ , and contains the inividual words of a sentence  $x$ , together with all possible labeled spans over that sentence:

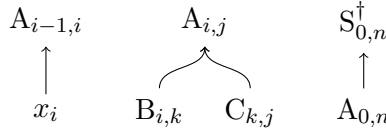
$$\mathcal{V} = \{x_i \mid 1 \leq i \leq n\} \cup \{(A, i, j) \mid A \in \Lambda, 0 \leq i < j \leq n\} \cup \{(S^\dagger, 0, n)\},$$

where  $(S^\dagger, 0, n)$  is a designated root node. The dependence on  $x$  can be made explicit by writing  $\mathcal{V}(x)$ .

The set of hyperedges  $\mathcal{E} \subseteq 2^\mathcal{V} \times \mathcal{V}$  specifies all the ways that adjacent constituents can be combined to form a larger constituent, under a particular grammar. Because we (implicitly) assume a normal form grammar that contains *all* possible productions, the set of hyperedges is particularly regular: the set  $\mathcal{E}$  contains all edges that connect nodes  $(B, i, k)$  and  $(C, k, j)$  at the tail with  $(A, i, j)$  at the head for all  $0 \leq i < k < j \leq n$ , all edges that connect  $x_i$  to  $(A, i, i + 1)$ , and all nodes can function as top nodes:

$$\begin{aligned} \mathcal{E} = & \left\{ \left\langle \{(B, i, k), (C, k, j)\}, (A, i, j) \right\rangle \mid A, B, C \in \Lambda, 0 \leq i < k < j \leq n \right\} \\ & \cup \left\{ \left\langle \{x_i\}, (A, i, i + 1) \right\rangle \mid A \in \Lambda, 1 \leq i \leq n \right\} \\ & \cup \left\{ \left\langle \{(A, 0, n)\}, (S^\dagger, 0, n) \right\rangle \mid A \in \Lambda \right\} \end{aligned}$$

The three kind of edges that make up  $\mathcal{E}$  are illustrated in figure 4.2. A tree is a set of nodes  $y \subseteq \mathcal{V}(x)$ , and the parse forest a set of trees  $\mathcal{Y}(x) \subseteq 2^{\mathcal{V}(x)}$ .



**Figure 4.2:** The three types of edges making up the hypergraph of example ??.

We connect a semiring  $\mathcal{K}$  to the hypergraph by defining the weight function as  $\omega : \mathcal{E} \rightarrow \mathbb{K}$ , and by accumulating the weights with its binary operations. The function  $\omega$  that assigns weights to the edges is given by either the function  $\psi$  or  $(\log \circ \psi)$ , depending on the semiring used. Because of this association, the function  $\omega$  has a very particular property: the function effectively depends only on the *head* of the edge. Given edges  $e = \langle \{u, w\}, v \rangle$  and  $e' = \langle \{u', w'\}, v \rangle$  for  $u \neq u'$  and  $w \neq w'$ , their weights are equal:  $\omega(e) = \omega(e')$ . For this reason we write  $\omega(v)$  instead.<sup>2</sup> This fact will allow us to greatly rewrite the recursion that follow.

With this structure in place, we are ready to derive the form of the inference algorithm particular to this structure.

### 4.3.2 Inside recursion

The inside recursion computes quantities  $\alpha(A, i, j)$  for all labels  $A \in \Lambda$  and all spans  $0 \leq i < j \leq n$ . The quantity computed depends on the semiring used. In this section we derive the inside recursion specific to our hypergraph from the general result given.

Let  $\mathcal{K}$  be some semiring with binary operations  $\oplus$  and  $\otimes$  and identity elements  $\hat{0}$  and  $\hat{1}$ . The inside recursion is given by the formula (Goodman,

---

<sup>2</sup>Thus implicitly defining  $\omega : \mathcal{V} \rightarrow \mathbb{K}$ .

1999)

$$\alpha(v) = \begin{cases} \hat{1} & \text{if } I(v) = \emptyset \\ \bigoplus_{e \in I(v)} \omega(e) \otimes \bigotimes_{u \in T(e)} \alpha(u) & \text{otherwise.} \end{cases}$$

At a node  $v = (A, i, i + 1)$  that spans one word  $x_i$ , the inside value is just the weight of the single edge incoming from that word:

$$\alpha(A, i, i + 1) = \omega(A, i, i + 1) \otimes \alpha(x_i) = \omega(A, i, i + 1), \quad (4.5)$$

for  $A \in \Lambda$ , for all  $0 \leq i < n$ . We used the fact that  $\alpha(x_i) = \hat{1}$ , which follows from the fact that there are no arrows incoming at  $x_i$ .

For a general node  $\alpha(A, i, j)$ ,  $j > i + 1$ , we observe that all the incoming edges have at the tail the nodes  $(B, i, k)$  and  $(C, k, j)$ , for all  $B, C \in \Lambda$  and  $i < k < j$ . The sum over edges thus reduces to independent sums over  $B$ ,  $C$ , and  $k$ , and the product over the inside values at the tail reduces to the product of values  $\alpha(B, i, k)$  and  $\alpha(C, k, j)$ . The form of  $\omega$  allows us to rewrite this greatly as

$$\begin{aligned} \alpha(A, i, j) &= \bigoplus_{B \in \Lambda} \bigoplus_{C \in \Lambda} \bigoplus_{k=i+1}^{j-1} \omega(A, i, j) \otimes \alpha(B, i, k) \otimes \alpha(C, k, j) \\ &= \omega(A, i, j) \otimes \bigoplus_{k=i+1}^{j-1} \bigoplus_{B \in \Lambda} \alpha(B, i, k) \otimes \bigoplus_{C \in \Lambda} \alpha(C, k, j) \\ &= \omega(A, i, j) \otimes \bigoplus_{k=i+1}^{j-1} \sigma(i, k) \otimes \sigma(k, j), \end{aligned} \quad (4.6)$$

where we've introduced the notational abbreviation

$$\sigma(i, j) = \bigoplus_{A \in \Lambda} \alpha(A, i, j).$$

Looking at 4.6 we can see the marginalized values  $\sigma(i, j)$  are all that are needed for the recursion. This suggest simplifying the recursion even further as

$$\begin{aligned} \sigma(i, j) &= \bigoplus_{A \in \Lambda} \alpha(A, i, j) \\ &= \left[ \bigoplus_{A \in \Lambda} \omega(A, i, j) \right] \otimes \left[ \bigoplus_{k=i+1}^{j-1} \sigma(i, k) \otimes \sigma(k, j) \right], \end{aligned} \quad (4.7)$$

where we put explicit brackets to emphasize that independence of the subproblems of labeling and splitting.

### 4.3.3 Outside recursion

The outside algorithm computes the quantities  $\beta(A, i, j)$  for all labels  $A \in \Lambda$  and all spans  $0 \leq i < j \leq n$ . For the labels spanning the entire sentence, the outside value is defined as

$$\beta(A, 0, n) = \begin{cases} \hat{1}, & \text{if } A = S^\dagger \\ \hat{0}, & \text{otherwise.} \end{cases}$$

To compute  $\beta(A, i, j)$  in the general case we need to sum over all possible upward binary expansions that  $A$  can be a part of, recursively making use of the outside values already computed. This means we sum over expansions  $(B \rightarrow C \ A, k, i, j)$  for all labels  $B$  and  $C$  and left endpoints  $1 \leq k < i$ , and sum over expansions  $(B \rightarrow A \ C, i, j, k)$  for all labels  $B$  and  $C$  and right endpoints  $j < k \leq n$ . This corresponds to the following recursive expression that uses the outside values of the superparts:

$$\begin{aligned} \beta(A, i, j) &= \bigoplus_{B \in \Lambda} \bigoplus_{C \in \Lambda} \bigoplus_{k=1}^{i-1} \psi(B, k, j) \otimes \alpha(C, k, i-1) \otimes \beta(B, k, j) \\ &\quad \oplus \bigoplus_{B \in \Lambda} \bigoplus_{C \in \Lambda} \bigoplus_{k=j+1}^n \psi(B, i, k) \otimes \beta(B, i, k) \otimes \alpha(C, j+1, k) \\ &= \bigoplus_{k=1}^{i-1} \left[ \bigoplus_{B \in \Lambda} \psi(B, k, j) \otimes \beta(B, k, j) \right] \otimes \left[ \bigoplus_{C \in \Lambda} \alpha(C, k, i-1) \right] \\ &\quad \oplus \bigoplus_{k=j+1}^n \left[ \bigoplus_{B \in \Lambda} \psi(B, i, k) \otimes \beta(B, i, k) \right] \otimes \left[ \bigoplus_{C \in \Lambda} \alpha(C, j+1, k) \right] \\ &= \bigoplus_{k=1}^{i-1} \sigma'(k, j) \otimes \sigma(k, i-1) \oplus \bigoplus_{k=j+1}^n \sigma'(i, k) \otimes \sigma(j+1, k) \end{aligned}$$

where

$$\begin{aligned}\sigma(i, j) &= \bigoplus_{A \in \Lambda} \alpha(A, i, j) \\ \sigma'(i, j) &= \bigoplus_{A \in \Lambda} \psi(A, i, j) \beta(A, i, j)\end{aligned}$$

#### 4.3.4 Solutions

Equipped with the two recursions and a handful of semirings we can provide the solutions promised at the outset of this section.

**Normalizer** When we instantiate the inside recursion with the real semiring, the value of  $\alpha$  at the root is the normalizer:

$$\alpha(S^\dagger, 0, n) = Z(x),$$

and when we instantiate the inside recursion with the log-real semiring we obtain the log-normalizer

$$\alpha(S^\dagger, 0, n) = \log Z(x).$$

**Parse** To find the viterbi tree  $\hat{y} = \arg \max_y p(y \mid x)$  and its probability  $p(\hat{y} \mid x)$  we use the Viterbi semirings (*cf.* examples C.9 and C.10 in appendix C). We take equation 4.7 and use the Viterbi semiring operations to derive that the value of the best subtree spanning words  $i$  to  $j$  is given by

$$\sigma(i, j) = \max_A [\log \psi(A, i, j)] + \max_k [\sigma(i, k) + \sigma(k, j)]. \quad (4.8)$$

The value  $\log \Psi(x, \hat{y})$  is then given by  $\sigma(0, n)$ , and can be normalized with to give the probability

$$p(\hat{y} \mid x) = \sigma(0, n) - \log Z(x). \quad (4.9)$$

The best label and splitpoint  $\hat{A}$  and  $\hat{k}$  for the span  $(i, j)$  are obtained by using the argmax:

$$\hat{A} = \arg \max_A \log \psi(A, i, j) \quad (4.10)$$

$$\hat{k} = \arg \max_k \sigma(i, k) + \sigma(k, j), \quad (4.11)$$

and the best tree  $\hat{y}$  is found by following back from the root down to the leaves the best splits and labels.

**Sample** Samples can be obtained by recursively sampling edges, starting at the root node  $(S^\dagger, 0, n)$ . The probability of an edge is proportional to the weight under that edge: this is precisely the inside value  $\alpha$  computed in the real-semiring. An edge  $e = \langle \{u, w\}, v \rangle$ , with  $u = (B, i, k)$ ,  $w = (C, k, j)$  and  $v = (A, i, j)$ , has probability

$$\begin{aligned} P(E = e) &= \frac{\omega(e) \otimes \bigotimes_{u \in T(e)} \alpha(u)}{\alpha(v)} \\ &= \frac{\psi(A, i, j) \alpha(B, i, k) \alpha(C, k, j)}{\alpha(A, i, j)} \\ &\propto \alpha(B, i, k) \alpha(C, k, j), \end{aligned} \tag{4.12}$$

where the edge-weight  $\psi(A, i, j)$  is absorbed into the normalizing constant because it is the same for all edges.

**Entropy** To compute the entropy  $H(Y \mid X = x)$  we need to first introduce the notion of the *marginal probability* of a node. The marginal of a node  $v = (A, i, j)$  in a hypergraph is the probability that it occurs in a tree  $y$  for the sentence  $x$  as governed by the distribution  $p$  that we defined on it. Let  $V$  be a random variable with the hypergraph nodes  $\mathcal{V}(x)$  as sample space, and define

$$\begin{aligned} P(V = v \mid X = x) &\triangleq \mathbb{E}_Y[\mathbf{1}_{\{v \in Y\}}] \\ &= \sum_{y \in \mathcal{Y}(x)} p(y \mid x) \mathbf{1}_{\{v \in y\}}. \end{aligned} \tag{4.13}$$

This is a probability distribution:

$$\begin{aligned} \sum_{v \in \mathcal{V}(x)} P(V = v \mid X = x) &= \sum_{y \in \mathcal{Y}(x)} p(y \mid x) \sum_{v \in \mathcal{V}(x)} \mathbf{1}_{\{v \in y\}} \\ &= \sum_{y \in \mathcal{Y}(x)} p(y \mid x) |y| \\ &\stackrel{?}{=} 1 \end{aligned}$$



The marginals can be computed from the inside and outside values computed in the real semiring as

$$P(V = v \mid X = x) = \frac{\alpha(A, i, j) \beta(A, i, j)}{Z(x)}, \quad (4.14)$$

a result from (Goodman, 1999).<sup>3</sup> The entropy can then be written as an expectation with respect to these marginals:

$$\begin{aligned} H(Y \mid X = x) &= - \sum_{y \in \mathcal{Y}(x)} p(y \mid x) \log p(y \mid x) \\ &= \log Z(x) - \sum_{y \in \mathcal{Y}(x)} p(y \mid x) \sum_{v \in y} \log \psi(x, v) \\ &= \log Z(x) - \sum_{y \in \mathcal{Y}(x)} p(y \mid x) \sum_{v \in \mathcal{V}(x)} \mathbf{1}_{\{v \in y\}} \log \psi(x, v) \\ &= \log Z(x) - \sum_{v \in \mathcal{V}(x)} \log \psi(x, v) \sum_{y \in \mathcal{Y}(x)} \mathbf{1}_{\{v \in y\}} p(y \mid x) \\ &= \log Z(x) - \sum_{v \in \mathcal{V}(x)} \log \psi(x, v) P(V = v \mid X = x) \\ &= \log Z(x) - \mathbb{E}_V[\log \psi(x, V)]. \end{aligned} \quad (4.15)$$

## 4.4 Experiments

We perform three types of experiments with the CRF parser:

- We show that the model is a good supervised parser. We train the model supervised on the PTB and show the f-score on the PTB test set.
- We evaluate the joint RNNG with samples from the CRF parser. We compare the perplexity and fscore with RNNG case.
- We evaluate ‘how good the model is’ as a sampler.

---

<sup>3</sup>This can also be seen by noting that the product of  $\alpha(A, i, j)$  and  $\beta(A, i, j)$  is the sum over all trees that contain the node  $v$ , and  $Z(x)$  the sum over all trees in general.

#### 4.4.1 Supervised model

We investigate the following.

- We have some optimization and hyperparameter choices here. The original paper uses Adam with 0.001 and a LSTM of dimension 250, which gives the model around 2.5 million parameters. For the discriminative RRNG we use SGD with 0.1, and hidden sizes of 128 gives the model around 800,000 parameters.
- I suggest two experiments: (1) use the default setting from (Stern *et al.*, 2017a) and (2) use the settings for the RRNG with a hidden size to match the 800,000 parameters.

#### 4.4.2 Proposal model

We investigate the following:

- We evaluate validation F-score and perplexity.
- We evaluate F-score with 100 samples (as many proposal trees as possible).
- We evaluate perplexity with varying number of samples: 1 (argmax), 10, 20, 50, 100 (default). The peplexity evaluation with the argmax prediction gives an impression of the uncertaty in the model (Buys and Blunsom, 2018).
- We perform learning rate decay and model selection based on a development score computed with the samples from the discriminative RRNG. Undecided: should we train a separate joint RRNG with CRF samples?

#### 4.4.3 Sampler

We investigate the following:

- We assess the conditional entropy of the model. This is most quantitative. Recall that conditional entropy is defined as

$$H(Y | X) = \sum_{x \in \mathcal{X}} p(x) H(Y | X = x), \quad (4.16)$$

where

$$H(Y|X = x) = - \sum_{y \in \mathcal{Y}} p(y|x) \log p(y|x). \quad (4.17)$$

The quantity  $H(Y|X = x)$  can be computed exactly with the CRF parser. We estimate the quantity  $H(Y|X)$  by a sum over the development dataset.

- We assess for some cherry picked sentences. This is more qualitative. These sentences should be difficult or ambiguous. Or they can be ungrammatical when taken from the syneval dataset. We can evaluate their entropy, and the diversity of samples, for example to see if there are clear modes. We can make violinplots of the probabilities of the samples. We can compute the f-scores of the samples compared with the argmax tree.

#### 4.4.4 Related work

The model that we presented regards a constituency tree as a collection of *labeled spans* over a sentence. Earlier CRF models for constituency parsing, both log-linear and neural, factorize trees over *anchored rules* (Finkel *et al.*, 2008; Durrett and Klein, 2015). This puts most of the expressiveness of the model in the state space of the dynamic program, modelling interactions between subparts of the trees through their interaction in the rules, instead of at the feature level. The model in Stern *et al.* (2017a) removes part of this structure, and puts more expressiveness in the input space by using rich neural feature representations conditioning on the entire sentence. The discrete interaction between the local scores remains only at level of labeled spans. This dramatically improves the speed of this model, which will become evident in the next section.

In recent decades, chart based parsing has moved from grammar to features: the work of Hall *et al.* (2014) shows how the log-linear CRF model of Finkel *et al.*, 2008 can work with bare unannotated grammars, when relying more on the surface features of the sentence, and Durrett and Klein (2015) show how these surface features can be replaced with dense neural features. The work of (Stern *et al.*, 2017a) moves one step further: the grammar is dispensed with altogether, making the model

span-factored, and the scoring function is given the full power of neural networks.

Contrast this with generative parsing based on treebank grammars, where features are not available because the models are not conditional. Instead, these models rely entirely on detailed rule information: basic treebank grammars do not parse well because the rules provide too little context, and good results can only be obtained by enriching grammars. The independence assumptions in the grammar are thus typically weakened, not strengthened. Such approaches lexicalize rules (Collins, 2003), annotate the rule with parent and sibling labels (Klein and Manning, 2003), or automatically learn refinements of nonterminal categories (Petrov *et al.*, 2006).

The closest predecessor to our model is the neural CRF parser of Durrett and Klein (2015), which predicts local potential for anchored rules using a feedforward network. This differs from our approach in two ways. Their method requires a grammar, extracted from a treebank beforehand, whereas our approach implicitly assumes all rules are possible rules in the grammar. Secondly, their scoring function conditions only on the parts of the sentence directly under the rule, dictated by the use of a feedforward network, whereas our scoring function computes a score based on representations computed from the entire sequence.

Earlier work on CRF parsing consider a tree as a collection of anchored rule productions productions Finkel *et al.*, 2008; Durrett and Klein, 2015, and hence define the score of a tree as the product over clique potentials on anchored rules:

$$\log \psi(A \rightarrow B \ C, i, k, j) = \log \psi(A, i, j) \quad (4.18)$$

$$(4.19)$$

discarding the rest of the span information. The function is then defined as

$$\log \psi(A, i, j) \triangleq s(i, j, A), \quad (4.20)$$

Note that the potential function as defined in ?? disregards most of the information in a binary rule. In particular we see that  $B$ ,  $C$  and  $k$ , the labels and split-point of the children, are discarded.

# 5

---

## Semisupervised learning

---

In this chapter we show how the RNNG can be trained on unlabeled data. Together with the regular, supervised, objective, this derives a way to perform semisupervised training.

- I formulate an unsupervised objective for the RNNG that we can combine with the supervised objective to perform semisupervised training.
- I introduce an approximate posterior in the form of a discriminative parser and derive a variational lower bound on the unsupervised objective.
- I show how to obtain gradients for this lowerbound by rewriting the gradient into a form that is called the score function estimator (Williams, 1992; Fu, 2006).

**Notation** In this chapter we write  $x$  for a sentence,  $y$  for a (latent) constituency tree, and  $\mathcal{Y}(x)$  for the *yield* of  $x$ , all trees that can be assigned to  $x$ . Furthermore, let  $\mathcal{D}_L = \{(x_i, y_i)\}_{i=1}^N$  be a labeled dataset of sentences  $x$  with gold trees  $y$ , and let  $\mathcal{D}_U = \{x_i\}_{i=1}^M$  be an unlabeled dataset consisting of just sentences  $x$ . We denote our generative RNNG

with  $p_\theta$  and when we write  $q_\lambda$  we will mean either the discriminative RNN or the CRF parser.

## 5.1 Objective

We define the following general semi-supervised objective

$$\mathcal{L}(\theta, \lambda) \triangleq \mathcal{L}_S(\theta) + \mathcal{L}_U(\theta, \lambda).$$

The supervised objective  $\mathcal{L}_S$  is optimized over  $\mathcal{D}_L$  and  $\mathcal{L}_U$  the unsupervised objective optimized over  $\mathcal{D}_U$ . We introduce  $\alpha \in \mathbb{R}_{\geq 0}$  as an arbitrary scalar controlling the contribution of the unsupervised objective.

**Supervised objective** We define the supervised objective  $\mathcal{L}_S(\theta)$  as

$$\mathcal{L}_S(\theta) \triangleq \sum_{(x,y) \in \mathcal{D}_L} \log p_\theta(x, y)$$

This objective is optimized as usual using stochastic gradient estimates:

$$\nabla_\theta \mathcal{L}_S(\theta) \approx \frac{N}{K} \sum_{(x,y) \in \mathcal{B}} \nabla_\theta \log p_\theta(x, y),$$

where  $\mathcal{B} \subseteq \mathcal{D}_U$  is a mini-batch of size  $K$  sampled uniformly from the dataset. We rely on automatic differentiation to compute  $\nabla_\theta \log p_\theta(x, y)$  (Baydin *et al.*, 2018).

**Unsupervised objective** We define the unsupervised objective  $\mathcal{L}_U(\theta, \lambda)$  as

$$\begin{aligned} \mathcal{L}_U(\theta, \lambda) &\triangleq \sum_{x \in \mathcal{D}_U} \log p(x) \\ &= \sum_{x \in \mathcal{D}_U} \log \sum_{y \in \mathcal{Y}(x)} p_\theta(x, y) \end{aligned}$$

This is a language modelling objective, in which we treat  $y$  as latent. We have noted the a consequence the lack of independence assumptions

of the RNNG is that the sum over trees  $y$  is not tractable. To optimize this objective we must thus fall back on approximate methods.

## 5.2 Variational approximation

We optimize the unsupervised objective using variational inference (Blei *et al.*, 2016). We introduce a posterior  $q_\lambda(y|x)$  parametrised by  $\lambda$  and use Jensen's inequality to derive a variational lower bound on the objective  $\mathcal{L}_U(\theta, \lambda)$ . First we bound the likelihood of one  $x$  in  $\mathcal{D}_U$

$$\begin{aligned} \log p(x) &= \log \sum_{y \in \mathcal{Y}(x)} q_\lambda(y|x) \frac{p_\theta(x, y)}{q_\lambda(y|x)} \\ &= \log \mathbb{E}_q \left[ \frac{p_\theta(x, y)}{q_\lambda(y|x)} \right] \\ &\geq \mathbb{E}_q \left[ \log \frac{p_\theta(x, y)}{q_\lambda(y|x)} \right] \\ &= \mathbb{E}_q \left[ \log p_\theta(x, y) - \log q_\lambda(y|x) \right] \end{aligned}$$

and write

$$\begin{aligned} \mathcal{E}(\theta, \lambda) &\triangleq \sum_{x \in \mathcal{D}_U} \mathbb{E}_q \left[ \log p_\theta(x, y) - \log q_\lambda(y|x) \right] \\ &\leq \sum_{x \in \mathcal{D}_U} \log p(x) \\ &= \mathcal{L}_U(\theta, \lambda) \end{aligned} \tag{5.1}$$

as a lower bound on our true objective  $\mathcal{L}_U$ . This is a particular instance of the evidence lower bound (ELBO) (Blei *et al.*, 2016). The quantity can be rewritten to reveal an entropy term  $H(q)$

$$\begin{aligned} \mathbb{E}_q \left[ \log p_\theta(x, y) - \log q_\lambda(y|x) \right] &= \mathbb{E}_q \left[ \log p_\theta(x, y) \right] - \mathbb{E}_q \left[ \log q_\lambda(y|x) \right] \\ &= \mathbb{E}_q \left[ \log p_\theta(x, y) \right] + H(q), \end{aligned}$$

which gives this objective an intuitive interpretation. On the one hand, the objective aims to

**Posterior** The posterior  $q$  can be any kind of models, with the only condition that for all  $x$  and  $y \in \mathcal{Y}(x)$ ,

$$p(x, y) > 0 \Rightarrow q(y|x) > 0.$$

This condition is fulfilled by any discriminatively trained parser with the same support as the joint RNNG  $p$ . We have two obvious choices at hand: the discriminatively trained RNNG, and the CRF parser that we introduced in chapter 4.

An interesting advantage of the CRF parser is that we *can* compute the entropy  $H(q)$  exactly. This contrasts with the discriminative RNNG, where  $H(q)$  can only be approximated. To make this explicit we introduce separate ELBO objectives:

$$\mathcal{E}_{\text{RNNG}}(\theta, \lambda) \triangleq \sum_{x \in \mathcal{D}_U} \mathbb{E}_q \left[ \log p_\theta(x, y) - \log q_\lambda(y|x) \right] \quad (5.2)$$

$$\mathcal{E}_{\text{CRF}}(\theta, \lambda) \triangleq \sum_{x \in \mathcal{D}_U} \mathbb{E}_q \left[ \log p_\theta(x, y) \right] + H(q). \quad (5.3)$$

### 5.3 Optimization

Just like the supervised objective  $\mathcal{L}_U$  we optimize the lower bound  $\mathcal{E}$  by gradient optimization, which means that we need to compute the gradients  $\nabla_\theta \mathcal{E}(\theta, \lambda)$  and  $\nabla_\lambda \mathcal{E}(\theta, \lambda)$ .

**Gradients of joint parameters** The first gradient is easy and permits a straightforward Monte-Carlo estimate:

$$\begin{aligned} \nabla_\theta \mathcal{E}(\theta, \lambda) &= \nabla_\theta \mathbb{E}_q \left[ \log p_\theta(x, y) - \log q_\lambda(y|x) \right] \\ &= \mathbb{E}_q \left[ \nabla_\theta \log p_\theta(x, y) \right] \\ &\approx \frac{1}{K} \sum_{i=1}^K \nabla_\theta \log p_\theta(x, y_i) \end{aligned}$$

where  $y_i \sim q_\lambda(\cdot|x)$  for  $i = 1, \dots, K$  are samples from the approximate posterior. We can move the gradient inside the expectation because  $q$  does not depend on  $\theta$ , and note that  $\nabla_\theta \log q_\lambda(y|x) = 0$ .



**Gradients of posterior parameters** The second gradient is not so straightforward and requires us to rewrite the objective into a form that is called the *score function estimator* (Fu, 2006). Firstly we define a *learning signal*

$$L(x, y) \triangleq \log p_\theta(x, y) - \log q_\lambda(y|x), \quad (5.4)$$

and use the identity in equation D.1 that we derive in the appendix

$$\begin{aligned} \nabla_\lambda \mathcal{E}(\theta, \lambda) &= \nabla_\lambda \mathbb{E}_q [L(x, y)] \\ &= \mathbb{E}_q [L(x, y) \nabla_\lambda \log q_\lambda(y|x)]. \end{aligned}$$

In this rewritten form the gradient is in the form of an expectation, and that does permit a straightforward MC estimate:

$$\mathbb{E}_q [L(x, y) \nabla_\lambda \log q_\lambda(y|x)] \approx \frac{1}{K} \sum_{i=1}^K L(x, y_i) \nabla_\lambda \log q_\lambda(x|y_i) \quad (5.5)$$

where again  $y_i \sim q_\lambda(\cdot|x)$  for  $i = 1, \dots, K$  are independently sampled from the approximate posterior. This estimator has been derived in slightly different forms in Williams (1992), Paisley *et al.* (2012), Mnih and Gregor (2014), Ranganath *et al.* (2014), and Miao and Blunsom (2016) and is also known as the REINFORCE estimator (Williams, 1992).

## 5.4 Variance reduction

We introduce the two baselines:

- Feedforward baseline (Miao and Blunsom, 2016).
- Argmax baseline from Rennie *et al.* (2017) which is exact in the CRF, and approximate in the RNN.
- The CRF has no variance in estimating the entropy.

## 5.5 Experiments

- Experiments with the two baselines and the two posteriors.

- Compare to a simple baseline: supervised learning on mixed gold-silver trees (partially predicted).
- Analyze the variance reduction provided by the different baselines.

## 5.6 Related work

- Discrete latent variables in neural models (Miao and Blunsom, 2016; Yin *et al.*, 2018).
- Semisupervised training for the RNNG (Cheng *et al.*, 2017).
- Argmax baseline Rennie *et al.*, 2017.
- Training with the policy gradient of a risk-objective as a surrogate for training with a dynamic oracle (Fried and Klein, 2018).

# 6

---

## Syntactic evaluation

---

Language models are typically evaluated by the perplexity they assign on held out data, and thus did we evaluate our language models in the previous chapters. In this chapter we look at alternatives. In particular, we look at evaluation that specifically probes the syntactic abilities of a language model. To this end we take the dataset introduced by Marvin and Linzen (2018), which presents a comprehensive set of syntactic challenges, and evaluate the models presented thus far against them. The dataset consists of constructed sentence pairs that differ in only one word, where one sentence is grammatical and the other is not, and the task is to assign higher probability to the grammatical sentence. We can think of this task as soliciting comparative *acceptability judgements*, which is a key concept in linguistics. We will refer to this dataset as Syneval, for *syntactic evaluation*.

**Organization** The chapter is organized as follows. First, I introduce the Syneval dataset and describe syntactic phenomena that it tests, and review how this dataset relates to other work on syntactic evaluation. I then describe multitask learning as an approach already suggested to improve language models for this task, and introduce a novel mul-

task model based on span labeling. Finally I evaluate all the models introduced in this thesis on this dataset and discuss the results.

## 6.1 Syntactic evaluation

A shortcoming of perplexity is that the metric conflates various sources of success (Marvin and Linzen, 2018). A language model can make use of many aspects of language to predict the probability of a sequence of words. And although we would like the probability of a sentence to depend on high-level phenomena such as syntactic well-formedness or global semantic coherence, a language model can also focus on lower hanging fruit such as collocations and other semantic relations predicted from local context. Syntax is especially difficult to evaluate given that most sentences in a corpus are grammatically simple (Marvin and Linzen, 2018). Arguably, this conflation is also the appeal: perplexity is a one size fits all metric. But for a fine-grained analysis we must resort to a fine-grained metric.

Recently, a series of papers has introduced tasks that specifically evaluate the syntactic abilities of language models (Linzen *et al.*, 2016; Gulordava *et al.*, 2018; Marvin and Linzen, 2018). And such tasks can be very revealing. The task introduced by Linzen *et al.* (2016) is to predict the correct conjugation of a verb given an earlier occurring subject—especially in the presence of distracting subjects that intervene<sup>1</sup>. This task has revealed that lower perplexity does not imply greater success on this task (Tran *et al.*, 2018), and that an explicitly syntactic model like the RNNG significantly outperforms purely sequential models, especially with increasing distance between the subject and the verb (Kuncoro *et al.*, 2018). This type of agreement is one of the phenomena evaluated in Syneval.

**Dataset** The Syneval dataset consists of contrastive sentence pairs that differ in only one word. Let  $(x, x')$  be this minimal pair, with grammatical sentence  $x$  and an ungrammatical sentence  $x'$ . Then a language model  $p$  makes a correct prediction on this pair if  $p(x) > p(x')$ .

---

<sup>1</sup>E.g. *Parts of the river valley have/has.*

The classification is based on the probability of the entire sequence. This makes the task applicable to grammatical phenomena that involve interaction between multiple words, or where the word of contention does not have any left context. This contrasts with the task introduced in Linzen *et al.* (2016). This approach is also more natural for a model like the RNNG, where the probability of the sentence is computed by marginalizing over all latent structures, whereas individual word probabilities can only be obtained when conditioning on a single structure—for example a predicted parse—as is done in Kuncoro *et al.*, 2018.

The sentence pairs fall into three categories that linguists consider to depend crucially on hierarchical syntactic structure (Everaert *et al.*, 2015; Xiang *et al.*, 2009):

1. Subject-verb agreement (The farmer *smiles*.)
2. Reflexive anaphora (The senators embarrassed *themselves*.)
3. Negative polarity items (*No* authors have *ever* been famous.)

The dataset contains constructions of increasing difficulty for each of these category. For example, the distance between two words in a syntactic dependency can be increased by separating them with a prepositional phrase: *The farmer next to the guards smiles*. In this example *the guards* additionally forms a distractor for the proper conjugation of *smiles*, making the example extra challenging.

The dataset is constructed automatically using handcrafted context-free grammars. The lexical rules are finegrained so that the resulting sentences reasonably coherent semantically. In particular there are rules for animate and inanimate objects so a sentence like *The apple laughs* cannot be constructed. The total dataset consists of around 350,000 sentence pairs.

**Categories** For the full list of categories that are evaluated in the Syneval dataset we refer the reader to appendix E.

**Related work** There has been a surge recently of work on syntactic evaluation. Linzen *et al.* (2016) introduce the task of long distance subject-verb agreement and Gulordava *et al.* (2018) make this test

more challenging by turning the sentences nonsensical while keeping them grammatical. Both datasets are extracted from a wikipedia corpus based on properties of their predicted dependency parse. To make the sentences nonsensical, Gulordava *et al.* (2018) randomly substitute words from the same grammatical category.<sup>2</sup> Warstadt *et al.* (2018) fine-tune neural models to learn to immitate grammatical acceptability judgments gathered from linguistics textbooks. McCoy *et al.* (2018) train a neural machine translation that turns a declarative sentence into a question, a kind of transformation that linguists have argued requires the existence of hierarchical structure in language (Everaert *et al.*, 2015).<sup>3</sup> Finally, targeted evaluation has been introduced previously for semantic comprehension by Zweig and Burges (2011) in a sentence completion task, and minimal contrastive sentence pairs have been used to evaluate neural machine translation by Sennrich (2017).

## 6.2 Multitask learning

One method that makes language models perform better in the tasks described in this chapter is to provide stronger syntactic supervision by using multitask learning (Enguehard *et al.*, 2017; Marvin and Linzen, 2018). Multitask learning is a simple yet effective way of providing additional supervision to neural network models by combining multiple tasks in a single objective while sharing parameters (Caruana, 1997), and has been succesfully applied in natural language processing (Collobert and Weston, 2008; Collobert *et al.*, 2011; Zhang and Weiss, 2016; Søgaard and Goldberg, 2016). By combining objectives that share parameters the model is encouraged to learn representations that are useful in all the tasks.

In this section I describe two simple baselines for the syntactic

---

<sup>2</sup>An approach inspired by Chomsky’s (in)famous sentence *Colorless green ideas sleep furiously* that is both grammatical and nonsensical.

<sup>3</sup>Such declarative-question pairs play a central role as empirical evidence in the argument—known as the *argument from the poverty of the stimulus*—that humans have an innate predisposition for generalizations that rely on hierarchical structure rather than linear order (Chomsky, 1980). Sequence-to-sequence neural machine translation, on the other hand, is a fully sequential model that involves no hierarchical structure or transformations.

evaluation task that are based on multitask learning. Both methods are based on language modelling with a syntactic side objective. The first side objective is to predict combinatory categorical grammar (CCG) supertags (Bangalore and Joshi, 1999) for each word in the sentence, and is proposed in (Enguehard *et al.*, 2017). The second side objective is to label spans of words with their category. This objective is inspired by the label scoring function in the CRF parser introduced in chapter 4. This is similar to an approach found in recent work on semantic parsing where a similar side-objective is used and where it is called a ‘syntactic scaffold’ (Swayamdipta *et al.*, 2018). The exact form of our side-objective is novel, as far as the author is aware of, and is parametrized in a considerably simpler way. Both objectives are challenging and should require representations that encode a fair amount of syntactic information.

**Multitask objective** In our case we combine a language model  $p$  with a syntactic model  $q$ , and optimize these jointly over a single labeled dataset  $\mathcal{D}$ <sup>4</sup>. In this case, we maximize the objective

$$\mathcal{L}(\theta, \lambda, \xi) = \sum_{(x,y) \in \mathcal{D}} \log p_{\theta, \lambda}(x) + \log q_{\theta, \xi}(y|x) \quad (6.1)$$

with respect to the parameters  $\theta$ ,  $\lambda$  and  $\xi$ . The key feature of multitask learning is that the two models  $p$  and  $q$  share the set of parameters  $\theta$  and that in objective 6.1 these parameters will thus be optimized to fit *both* the models well. The parameters in  $\lambda$  and  $\xi$  are optimized for their respective objectives separately. The proportion and the nature of the parameters that belong to  $\theta$  is a choice of the modeller and determines how both objectives influence one another. In the following paragraphs we specify this parametrization.

---

<sup>4</sup>Note that it is our choice to focus on only a single dataset, and that multitask learning is principle more flexible than that, providing the option to combine multiple disjoint datasets in a single objective.

**Language model** The main model  $p$  is a regular RNN language model on sentences  $x$ :

$$\log p_{\theta,\lambda}(x) = \sum_{i=1}^n \log p_{\theta,\lambda}(x_i \mid x_{<i}),$$

where the conditional probabilities of  $x_i$  are computed by linear regression on forward RNN vectors  $\mathbf{f}_{i-1}$  as

$$p_{\theta,\lambda}(x \mid x_{<i}) \propto \exp \left[ \mathbf{W}^\top \mathbf{f}_{i-1} + \mathbf{b} \right]_x.$$

The parameters  $\xi = \{\mathbf{W}, \mathbf{b}\}$  are specific to the language model, and the vectors  $\mathbf{f}_i$  are computed using a forward RNN parametrized by  $\theta$ . The vectors  $\mathbf{f}$  are used in the side objective as well, and it is in this precise sense that the parameters  $\theta$  are shared between  $p$  and  $q$ .

**Word labeling** Let  $y = (y_1, \dots, y_n)$  be a sequence of CCG supertags for the sentence  $x$ , with one tag for each word. The side model is then a simple greedy tagging model:

$$\log q_{\theta,\xi}(y \mid x) = \sum_{i=1}^n \log q_{\theta,\xi}(y_i \mid x).$$

The probability over tags for position  $i$  are computed from  $\mathbf{f}_i$  using a feedforward network parametrized by  $\xi$

$$q_{\theta,\xi}(y_i \mid x) \propto \exp \left[ \text{FFN}_\xi(\mathbf{f}_i) \right]_{y_i}.$$

This side objective is taken from Enguehard *et al.* (2017) and is the side objective that is used in the multitask model of Marvin and Linzen (2018).

**Span labeling** Let  $y$  be a sequence of labeled spans  $y_k = (\ell_k, i_k, j_k)$  that are obtained from a gold parse tree for  $x$ . Given span endpoints  $i$



and  $j$  and the sentence  $x$ , the side model  $q$  predicts a label  $\ell$ :

$$\begin{aligned}\log q_{\theta,\xi}(y \mid x) &= \sum_{k=1}^K \log q_{\theta,\xi}(y_k \mid x) \\ &= \sum_{k=1}^K \log q_{\theta,\xi}(\ell_k \mid x, i_k, j_k).\end{aligned}$$

The representation for the span is defined as in the CRF parser as

$$\mathbf{s}_{ij} = \mathbf{f}_j - \mathbf{f}_i,$$

and the distribution over labels is computed using a feedforward network parametrized by  $\xi$  as

$$q_{\theta,\xi}(\ell \mid x, i, j) \propto \exp[\text{FFN}_{\xi}(\mathbf{s}_{ij})]_{\ell}.$$

This model differs from the supertagging model in the interaction of the vectors  $\mathbf{f}$  in the definition of  $vecs_{ij}$ . This linear definition puts significant restriction on the encodings, which is precisely what we want. This side objective is inspired by the label scoring function of the CRF parser, and is novel as a multitask objective, as far as I know.

## 6.3 Experiments

In this section we report the results on the dataset.

### 6.3.1 Baselines

Marvin and Linzen (2018) provide a number of baselines:

- Human evaluation via amazon mechanical turk.
- N-gram language model with kneser ney smoothing.
- RNN language model (LSTM).
- RNN language model with multitask learning.

We reproduce the RNN language model and the RNN multitask language model,

### 6.3.2 Setup

- Following the example of Marvin and Linzen (2018) we train a regular LSTM language model, and a multitask LSTM language model with CCG supertagging. We additionally train
- For the language train 10 independent runs of each model and report means and standard deviations of their accuracy on the Syneval dataset.

Details about training can be found in appendix [B](#)

### 6.3.3 Results

# 7

---

## Conclusion

---

Here is a narrative summary of what I have shown in this thesis.

### 7.1 Main contributions

The main  $n$  contributions of thesis are:

**Global training of a chart based neural parser.** Here I describe what that entails.

**Semisupervised training of RNNGs.** Here I describe what that entails.

**Effective baselines for the score function estimator.** Here I describe what that entails.

### 7.2 Future work

We have identified possibilities for future work:

**Something.** Here I describe what that entails.

## **Appendices**

# A

---

## Figures

---

In this appendix I will put figures, for cases where there are just too many. For example:

- The barplots of the syntactic evaluation
- The training losses for the various models
- The valuation perplexity and f-score during training.

### A.1 Training plots

We show training plots that are means with standard deviation bands over 10 runs. We have two types of plots: losses and development scores.

**Development scores** Plots with development scores.

- DiscRNNG + GenRNNG-disc + GenRNNG-crf + CRF (small) development fscores.
- CRF parser losses: our (128d + SGD) and original (250d + Adam).

- GenRNNG-disc + GenRNNG-crf development perplexity
- Language models: lstm + lstm-ccg + lstm-span development perplexity
- GenRNNG + LSTM development perplexity.

**Training losses** Plots with training losses.

## **A.2 Test scores**

Here we show violin plots that show the distribution of the test scores.

## **A.3 Samples**

Here we put figures to illustrate the samples.

## **A.4 Syntactic evaluation**

Here we put more bar-charts for the syntactic evaluation.

# B

---

## Implementation

---

This appendix reports all details about the implementation of the models, including the datasets and their pre-processing, optimization, and hyperparameters. We have made some deliberate choices which we will motivate. As a general rule, we choose simplicity over maximal performance: we use a minimal scheme for dealing with unknown words, we use only word embeddings, which we learn from scratch, and use regular SGD optimization. The reason: to compare the models in a minimal setting so as to maximally focus on the essential modelling differences between the models. All code is available at [github.com/daandouwe/thesis](https://github.com/daandouwe/thesis).

### B.1 Data

#### B.1.1 Datasets

**Penn Treebank** We use a publicly available version of the Penn Treebank (PTB) that was preprocessed and published by Cross and Huang (2016), and that has since been used in the experiments of Stern *et al.* (2017a) and Kitaev and Klein (2018). The data is divided along the standard splits of sections 2-21 for training, section 22 for development,

and section 23 for testing. The dataset comes with predicted tags, which is a requirement for neural parsers to avoid overfitting, but note however that none of our models use tag information. The dataset is available at <https://github.com/jhcross/span-parser/data>.<sup>1</sup>

**One Billion Word Benchmark** In the experiments on semisupervised learning we use the One Billion Word Benchmark (OBW) dataset (Chelba *et al.*, 2013) to obtain unlabeled data. This is a common dataset for large-scale language modelling (Jozefowicz *et al.*, 2016), and has the advantage that it has sentences separated by a newline; a requirement for the RNN language model, which should only be applied to entire sentences and not to longer or shorter segments.<sup>2</sup> The dataset in its entirety is far too large for our purposes, so instead we select sentences from the first section of the training data<sup>3</sup> by selecting the first 100,000 sentences that have at most 40 words. Because the OBW uses slightly different tokenization and standardization of characters than the Penn Treebank we perform a number of processing steps to smooth out these differences. First, we escape all brackets following the Penn Treebank convention (*replacing (with -LRB-) and do the same for the quotation marks (replacing " with ").* Finally, tokenization of negation is handled differently in the OBW, and we change this to the PTB convention (*replacing don't with do n't*). These simple changes together avoid a lot of incoherence and scripts for preprocessing are available at [github.com/daandouwe/thesis/scripts](https://github.com/daandouwe/thesis/scripts).

**CCG supertags** For the multitask language model with CCG supertagging we use the CCGBank (Hockenmaier and Steedman, 2007) processed by Enguehard *et al.* (2017) into a word-tag format. This is the dataset also used by Marvin and Linzen (2018). It follows the same splits of the Penn Treebank as described above, and restricts the size of the tagset

---

<sup>1</sup>Although I do not know how it is possible to make the PTB public, given the licensing restrictions of the LDC, I am very thankful that it was done. Now, all the data used in this thesis is publicly available.

<sup>2</sup>For this reason we cannot use the otherwise appealing Wikitext dataset (Merity *et al.*, 2016): this dataset has sentences grouped into paragraphs.

<sup>3</sup>Section `news.en-00001-of-00100`.



from the original 1363 different supertags to 452 supertags that occurred at least ten times, replacing the rest of the tags with a dummy token. The dataset is publicly available at [https://github.com/BeckyMarvin/LM\\_syneval/tree/master/data/ccg\\_data](https://github.com/BeckyMarvin/LM_syneval/tree/master/data/ccg_data).

### B.1.2 Vocabularies

We use two types of vocabularies corresponding to the two types of models that we study in this thesis: a vocabulary for the discriminative models and a vocabulary for the generative models. The vocabulary used in the discriminative models contains all words in the training data, whereas the vocabulary used in the generative models only includes words that occur at least 2 times in the training data<sup>4</sup>. Finally, in the semisupervised models we construct the vocabulary from the labeled and unlabeled datasets combined. To keep the vocabulary size manageable in this larger dataset we restrict the vocabulary of the generative model to words that occur at least 3 times.

**Unknown words** We use a single token for unknown words, and during training replace each word  $w$  by this token with probability

$$\frac{1}{1 + \text{freq}(w)},$$

where  $\text{freq}(w)$  is the frequency of  $w$  in the training data. In this we follow Stern *et al.* (2017a) and deviate from Dyer *et al.* (2016), who use a set of almost 50 tokens each with detailed lexical information about the unknown word in question.<sup>5</sup> This elaborate approach is common in parsing but certainly not in language modelling (Dyer *et al.*, 2016), for which reason we opt for the simpler scheme of a single token.

**Embeddings** All word embeddings are learned from scratch: the embeddings are considered part of the model’s parameters and are optimized jointly with the rest of them, starting from random initialization

---

<sup>4</sup>This makes training of the generative model faster, because the softmax normalization involves less terms in the sum, and additionally avoids the statistical difficulty related to predicting words that occur just once. The discriminative model has neither of these problems, since the words are only conditioned on

<sup>5</sup>An approach taken from the Berkeley parser (Petrov *et al.*, 2006).

Table B.1: Vocabularies

(Glorot and Bengio, 2010). We surmise that more elaborate embeddings should improve performance of the models,<sup>6</sup> but such investigation is in principle orthogonal to our work. We furthermore do not experiment with any kind of subword information in the embeddings and, as noted before, make no use of tags in any of the models. The influence of such embeddings on the discriminate parser of Stern *et al.* (2017a) is analysed extensively by Gaddy *et al.* (2018), who investigate all combinations of word, tag, and character-LSTM embeddings, and find that the best model<sup>7</sup> improves on the worst model<sup>8</sup> by only 0.8 F1, which is a relative improvement of just 1%. We believe this justifies our basic approach.

## B.2 Implementation

All our models are implemented in python using the Dynet neural network library (Neubig *et al.*, 2017b), and use automatic batching (Neubig *et al.*, 2017a). Autobatching enables efficient training of our models, for which manual batching is not possible.

**Optimization** All our models are optimized with stochastic gradient-based methods, in which we use mini-batches to compute stochastic approximations of the model’s gradient on the entire dataset. We use mini-batches subsampled from the dataset and let Dynet compute the gradients using automatic differentiation (Neubig *et al.*, 2017b; Baydin *et al.*, 2018). For all our supervised models we use regular stochastic gradient descent (SGD), with an initial learning rate of 0.1, and anneal this based on performance on a development set. This follows the recommendations of Wilson *et al.* (2017), who show that on models and datasets similar to ours this method finds good solutions and is more

---

<sup>6</sup>See for example the impact of ELMo embeddings (Peters *et al.*, 2018) on the performance of the parser in Kitaev and Klein (2018) (an adaptation of the chart parser in Stern *et al.* (2017a)).

<sup>7</sup>A tie between the model that uses all embeddings concatenated and the model that uses just the character-LSTM, both with an F1 of 92.24.

<sup>8</sup>Using only word embeddings, at an F1 of 91.44.

**Table B.2:** Optimization**Table B.3:** Model parameters

robust against overfitting than methods with adaptive learning rates. This also follows Dyer *et al.* (2016), who obtain their best RNNG models using this optimizer and learning rate. For our semisupervised model we do rely on adaptive gradient methods, and use Adam (Kingma and Ba, 2014) with the default learning rate<sup>9</sup> of 0.001. Adaptive learning are considered more suitable for dealing with the dramatic variability in magnitude of the surrogate objective (Ranganath *et al.*, 2014; Fried and Klein, 2018). We use minibatches

We normalize the learning signal: “This variability makes training an inference network using a fixed learning rate difficult. We address this issue by dividing the centered learning signal by a running estimate of its standard deviation. This normalization ensures that the signal is approximately unit variance, and can be seen as a simple and efficient way of adapting the learning rate.” (Mnih and Gregor, 2014).

Surrogate objective and gradient blocking Schulman *et al.*, 2015.

**Hyperparameters** For the RNNG we follow exactly the hyperparameters from the published papers. For the CRF parser

- Dropout of ...

---

<sup>9</sup>To be precise, this is the value  $\alpha$  in Kingma and Ba (2014).

# C

---

## Semiring parsing

---

The derivation of the inference algorithms in chapter 4 makes use of the notion of a *weighted hypergraph* as a compact representation of all parses and their scores (Gallo *et al.*, 1993; Klein and Manning, 2004), and also makes use some of the ideas and notation of *semiring parsing* (Goodman, 1999; Li and Eisner, 2009).

### C.1 Hypergraph

We use a *backward hypergraph* to compactly represent all possible trees over a sentence under a grammar, an idea introduced in Klein and Manning (2004). We call such hypergraph a parse forest.

**Definition C.1.** A *directed hypergraph* is a pair  $G = (V, E)$ , consisting of a set of nodes  $V$ , and a set of directed hyperedges  $E \subseteq 2^V \times 2^V$  that connect a set of nodes at the *tail* of the arrow to a set of node at the *head* of the arrow.

**Definition C.2.** A *backward-hypergraph* is a directed hypergraph where the hyperedges  $E \subseteq 2^V \times V$  connect to a single node. For a node  $v \in V$  the set  $I(v) \subseteq E$  denotes the set of edges incoming at  $v$  and the set  $O(v) \subseteq E$  denotes the set of edges outgoing from  $v$ , *i.e.* the edges for

which  $v$  is in the tail. For an edge  $e$  we write  $T(e) \subseteq V$  for the set of nodes in the tail of  $e$ , and define  $H(e) \in V$  as the node at the head of the edge.

**Definition C.3.** A *weighted hypergraph* is any hypergraph  $G$  equipped with a weight function  $\omega : E \rightarrow K$  that to each edge assigns a weight from a weight-set  $K$ . Typically  $K$  the set of real numbers, or the set of integers.

**Definition C.4.** A *hyperpath* in a backward-hypergraph is a set of edges that connects a single node  $v \in V$  at the sink to a set of nodes  $W \subseteq V$  upstream. Formally defining a hyperpath is a little cumbersome, but informally, a hyperpath is a set of edges obtained by the following recursive procedure: starting at  $v$ , follow a single incoming hyperarc  $e$  backwards, save the edge and repeat this for all nodes in  $T(e)$ ; stop when all nodes in  $W$  have been encountered. In the context of parsing we also call a hyperpath from the root node to the set of words  $\{x_1, \dots, x_n\}$  a *derivation*: this hyperpath represents a single tree for the sentence.

**Example C.1.** (Parse forest) Figure C.1 shows a fragment of a hypergraph that represents the parse forest over an example sentence. Shown are the two hyperpaths corresponding to the two partially overlapping derivations

$$(S \text{ (NP } The \text{ (ADJP } very \text{ hungry) } cat) \text{ (VP } meows) \text{ )}) \quad (C.1)$$

$$(S \text{ (NP } The \text{ (NP (ADJP } very \text{ hungry) } cat)) \text{ (VP } meows) \text{ )}) \quad (C.2)$$

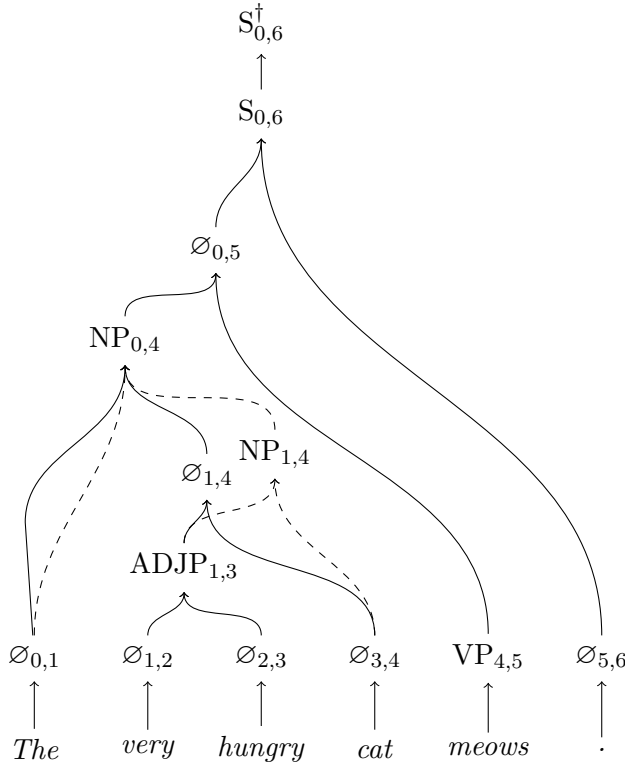
after collapsing the empty nodes  $\emptyset$ . The edges that direct to the special node  $(S^\dagger, 0, n)$  make the hypergraph *rooted*: each *hyperpath* with the set of all words at the source ends at this one node.

## C.2 Semiring

We use *semirings* to compute various quantities of interest over a weighted hypergraph.

**Definition C.5.** A *semiring* is an algebraic structure

$$\mathcal{K} = (\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1}),$$



**Figure C.1:** A fraction of a parse hypergraph showing two possible parses.

over a field  $\mathbb{K}$ , with additive and multiplicative operations  $\oplus$  and  $\otimes$ , and additive and multiplicative identities  $\bar{0}$  and  $\bar{1}$ . A semiring is equivalent to a ring<sup>1</sup>, without the requirement of an additive inverse for each element.

Some of the semirings relevant to our discussion are the following:

**Example C.2.** The *real semiring*

$$(\mathbb{R}_{\geq 0}, +, \times, 0, 1),$$

defined over the nonnegative reals, with regular addition and multiplication is a semiring. Note that the additive inverse is missing for all

<sup>1</sup>Perhaps the most common algebraic structure around: the real numbers with regular addition and multiplication form a ring.

elements greater than zero (*i.e.* missing the negative reals).

**Example C.3.** The *boolean semiring*

$$(\{\top, \perp\}, \vee, \wedge, \top, \perp),$$

is defined over truth values  $\top$  (True, or 1), and  $\perp$  (False, or 0). The binary operations are the logical ‘and’ and ‘or’ operations.

**Example C.4.** The *log-real semiring*,

$$(\mathbb{R} \cup \{-\infty\}, \oplus, +, -\infty, 0),$$

defined over the reals extended including  $\infty$ , with addition defined as the logarithmic sum<sup>2</sup>

$$a \oplus b = \log(e^a + e^b),$$

and multiplication is defined as regular addition.

**Example C.5.** The *max-tropical semiring*, or *Viterbi semiring*<sup>3</sup>

$$(\mathbb{R}_{\geq 0} \cup \{-\infty\}, \max, +, -\infty, 0),$$

is the log-real semiring that uses the max operator for addition.

### C.3 Semiring parsing

A semiring  $\mathcal{K}$  can be connected to a weighted hypergraph by defining the function  $\omega$  over its field  $\mathbb{K}$ , and by accumulating the weights with its binary operations. When the hypergraph represents a parse forest, we are in the realm of *semiring parsing* (Goodman, 1999).

The key result derived by Goodman (1999) is that many quantities of interest can be computed by a single recursion but over different semirings. First, we establish the relation between the weight function and the structures encoded in the hypergraph by defining the weight of a derivation and the weight of an entire hypergraph.

---

<sup>2</sup>Also known as *log-sum-exp*, or *log-add-exp*.

<sup>3</sup>This naming will become clear.

**Definition C.6.** (Hypergraph weights) Let  $G_\omega = (V, E, \omega)$  be a weighted hypergraph, with  $\omega$  defined over a semiring  $\mathcal{K}$ . We define the weight of the derrivation  $D \subseteq E$  as the product of the weights of the edges:

$$\bigotimes_{e \in D} \omega(e). \quad (\text{C.3})$$

Let  $\mathcal{D} \subseteq 2^E$  be the set of all derivations in the hypergraph  $G$ . Then the total weight of the hypergraph under  $\omega$  is defined as the sum of the weights of all the derrivations in it:

$$\bigoplus_{D \in \mathcal{D}} \bigotimes_{e \in D} \omega(e). \quad (\text{C.4})$$

**Example C.6.** (CRF parser) The CRF parser assigns a score  $\Psi(x, y) \geq 0$  by factorizing over parts of  $y = \{y_a\}_{a=1}^A$  as a product of potentials  $\psi(x, y_a) \geq 0$  of the parts. The parts  $y_a$  are the edges in the hyperpath that create the derivation  $y$ . In our CRF, thus, the weight function  $\omega$  is given by the function  $\psi$ ; the function  $\Psi$  is equivalent to equation C.3; and the sum over  $\mathcal{Y}(x)$  is equivalent to equation C.4.

### C.3.1 Inside and outside recursions

We follow the exposition of Li and Eisner (2009), but the results where first derived in Goodman (1999).

**Definition C.7.** The *inside value*  $\alpha(v)$  at a node  $v \in V$  accumulates the weight of all the paths that converge at that node. The accumulation is relative to a semiring, and is defined as

$$\alpha(v) = \begin{cases} \hat{1} & \text{if } I(v) = \emptyset \\ \bigoplus_{e \in I(v)} \omega(e) \otimes \bigotimes_{u \in T(e)} \alpha(u) & \text{otherwise.} \end{cases}$$

The value  $\alpha(v)$  at the root node  $v$  is the sum of the weight of all the derivations in the hypergraph. The recursion can solved by visiting the nodes in  $V$  in topological order.

**Definition C.8.** The *outside value*  $\beta(v)$  at a node  $v \in V$  accumulates over the weights of all the paths that head out from  $v$ . The accumulation



is relative to a semiring, and is defined as:

$$\beta(v) = \begin{cases} \hat{1} & \text{if } O(v) = \emptyset \\ \bigoplus_{e \in O(v)} \omega(w) \otimes \beta(H(e)) \otimes \bigotimes_{\substack{w \in T(e) \\ w \neq u}} \alpha(w) & \text{otherwise.} \end{cases}$$

The recursion can be solved by visiting the nodes in  $V$  in reverse topological order.

### C.3.2 Instantiated recursions

By instantiating the inside and outside recursions with different semirings we solve a whole range of problems.

**Example C.7.** (Inside and outside values) When we instantiate the inside recursion and the outside recursion with the *real semiring*, the algorithms reduce to the classical inside-outside algorithm (Baker, 1979). The values  $\alpha(v)$  and  $\beta(v)$  are the inside and outside values. In particular, the value of  $\alpha$  at the root is the normalizer:

$$\alpha(S^\dagger, 0, n) = Z(x).$$

In our CRF parser, the function  $\omega$  is given by the nonnegative function  $\Psi$ .

**Example C.8.** (Logarithmic domain) When we are concerned with numerical stability, or we only need the logarithm of the quantities of interest, we can use the *log-real semiring*. The values  $\alpha(v)$  and  $\beta(v)$  now give the log-inside and log-outside values respectively. In particular, the value  $\alpha$  at the root now gives the lognormalizer:

$$\alpha(S^\dagger, 0, n) = \log Z(x).$$

The semiring addition  $ab = \log(e^a + e^b)$  is made numerically stable by writing

$$\log(e^a + e^b) = a + \log(1 + e^{b-a}) = b + \log(1 + e^{a-b})$$

and choosing the expression with the smaller exponent. In the CRF, the function  $\omega$  is given by the composed function  $\log \circ \Psi$ .

**Example C.9.** (Viterbi weight) When we instantiate the inside recursion with the *max-tropical semiring* we get the recursion that computes at each node the subtree of maximum weight. The value  $\alpha$  at the root is the weight of the derivation with the maximum weight

$$\alpha(S^\dagger, 0, n) = \log \Psi(x, \hat{y}),$$

where  $\hat{y}$  is the Viterbi tree. Normalizing the score with the lognormalizer gives the log-probability

$$\log p(\hat{y} \mid x) = \alpha(S^\dagger, 0, n) - \log Z(x).$$

Hence the alternative name *Viterbi semiring*.

**Example C.10.** (Viterbi derivation) The Viterbi semiring derives the *weight* of the best tree. Replacing the max in the Viterbi semiring with an argmax derives the best tree itself:

$$\alpha(S^\dagger, 0, n) = \hat{y}, \quad \hat{y} \triangleq \arg \max_{y \in \mathcal{Y}(x)} p(y \mid x).$$

Roughly speaking, because although this idea can be made precise by constructing the *Viterbi-derivation semiring* (Goodman, 1999) over the set of possible derivations, with corresponding binary operations and set-typed identity elements  $\hat{0}$  and  $\hat{1}$ , it is a little more complicated than that<sup>4</sup>.

**Example C.11.** (Recognition) When we instantiate the inside recursion with the *boolean semiring* we get the recursion that recognizes whether a hyperpath exists from the words spanned by the node:

$$\alpha(S^\dagger, 0, n) = \begin{cases} \top & \text{if } x \in L(G), \\ \perp & \text{otherwise.} \end{cases}$$

The value at the root tells whether the sentence has at least one parse. Due to the trivial grammar used this will always be  $\top$  in our CRF.

---

<sup>4</sup>And most of all, a bit cumbersome.

# D

---

## Variational inference

---

In this appendix we give an account of variational inference in general of amortized variational inference in particular. We focus on amortized inference with discrete latent variables, and in particular when the variables are structured. We then derive the score function gradient, as used in chapter 5, and we describe techniques to reduce the variance of this estimator.

### D.1 Variational Inference

We will write some generic things about variational inference.

- Variational inference for exponential families, with conjugate priors (Jordan *et al.*, 1999; Wainwright and Jordan, 2008; Blei *et al.*, 2016).
- With amortized inference, using neural networks, for non-conjugate models (Kingma and Welling, 2014; Rezende *et al.*, 2014) and in particular the reparametrization trick that makes these models efficiently trainable.
- Reparametrization for discrete latent variables (Maddison *et al.*,

2017; Jang *et al.*, 2017).

- The generalization of this reparametrization trick in automatic differentiation variational inference (Kucukelbir *et al.*, 2017).
- Black box variational inference, which uses the same combination of score function gradient with variance reduction that we resort to (Ranganath *et al.*, 2014).

## D.2 Score function estimator

In this section we provide a detailed derivation of the score function estimator:

$$\nabla_{\lambda} \mathbb{E}_q [L(x, y)] = \mathbb{E}_q [L(x, y) \nabla_{\lambda} \log q_{\lambda}(y|x)] \quad (\text{D.1})$$

where

$$L(x, y) \triangleq \log p_{\theta}(x, y) - \log q_{\lambda}(y|x)$$

$$\begin{aligned} \nabla_{\lambda} \mathbb{E}_q [L(x, y)] &= \nabla_{\lambda} \mathbb{E}_q [\log p_{\theta}(x, y) - \log q_{\lambda}(y|x)] \\ &= \nabla_{\lambda} \sum_{y \in \mathcal{Y}(x)} \left\{ q_{\lambda}(y|x) \log p_{\theta}(x, y) - q_{\lambda}(y|x) \log q_{\lambda}(y|x) \right\} \\ &= \sum_{y \in \mathcal{Y}(x)} \left\{ \nabla_{\lambda} q_{\lambda}(y|x) \log p_{\theta}(x, y) - \nabla_{\lambda} q_{\lambda}(y|x) \log q_{\lambda}(y|x) \right. \\ &\quad \left. - q_{\lambda}(y|x) \nabla_{\lambda} \log q_{\lambda}(y|x) \right\} \\ &= \sum_{y \in \mathcal{Y}(x)} \left\{ \nabla_{\lambda} q_{\lambda}(y|x) \log p_{\theta}(x, y) - \nabla_{\lambda} q_{\lambda}(y|x) \log q_{\lambda}(y|x) \right\} \\ &= \sum_{y \in \mathcal{Y}(x)} \left\{ L(x, y) \nabla_{\lambda} q_{\lambda}(y|x) \right\} \\ &= \sum_{y \in \mathcal{Y}(x)} \left\{ L(x, y) q_{\lambda}(y|x) \nabla_{\lambda} \log q_{\lambda}(y|x) \right\} \\ &= \mathbb{E}_q [L(x, y) \nabla_{\lambda} \log q_{\lambda}(y|x)]. \end{aligned}$$

In this derivation we used the identity

$$\nabla_{\lambda} q_{\lambda}(y|x) = q_{\lambda}(y|x) \nabla_{\lambda} \log q_{\lambda}(y|x),$$

which follows from the derivative

$$\nabla_{\lambda} \log q_{\lambda}(y|x) = \nabla_{\lambda} q_{\lambda}(y|x) q_{\lambda}(y|x)^{-1}.$$

We

$$\begin{aligned} \sum_{y \in \mathcal{Y}(x)} q_{\lambda}(y|x) \nabla_{\lambda} \log q_{\lambda}(y|x) &= \sum_{y \in \mathcal{Y}(x)} q_{\lambda}(y|x) \frac{\nabla_{\lambda} q_{\lambda}(y|x)}{q_{\lambda}(y|x)} \\ &= \sum_{y \in \mathcal{Y}(x)} \nabla_{\lambda} q_{\lambda}(y|x) \\ &= \nabla_{\lambda} \sum_{y \in \mathcal{Y}(x)} q_{\lambda}(y|x) \\ &= \nabla_{\lambda} 1 \\ &= 0. \end{aligned}$$

### D.3 Optimization

We use automatic differentiation ([baydin2017automatic](#)) to obtain all our gradients. In order to obtain the gradients in formula [D.1](#) using this method we rewrite it in the form of a *surrogate objective* (Schulman *et al.*, [2015](#)):

$$\mathcal{L}_{\text{Surr}}(\theta, \lambda) = \frac{1}{K} \sum_{i=1}^K \log q_{\lambda}(x|y_i) \text{BLOCKGRAD}(L(x, y_i)). \quad (\text{D.2})$$

The function `BLOCKGRAD` detaches a node from its upstream computation graph. This turns it effectively into a scalar. More precisely, let  $f$  be function (computed by a node in the computation graph) with parameters  $\theta$  and input  $x$ , then

$$\text{BLOCKGRAD}(f_{\theta}(x)) \triangleq f(x),$$

such that

$$\nabla_{\theta} \text{BLOCKGRAD}(f_{\theta}(x)) = \nabla_{\theta} f(x) = 0.$$

Automatic differentiation of equation D.2 with respect to  $\lambda$  will give us the exact expression we are looking for

$$\nabla_{\lambda} \mathcal{L}_{\text{SURR}}(\theta, \lambda) = \frac{1}{K} \sum_{i=1}^K L(x, y) \nabla_{\lambda} \log q_{\lambda}(x|y_i),$$

hence the adjective *surrogate*.

## D.4 Variance reduction

We have derived an estimator for the gradient of the posterior parameters in the unsupervised objective. This estimator is unbiased, but is known to have high variance, often too much to be useful (Paisley *et al.*, 2012). Two effective methods to counter this are control variates and baselines (Ross, 2006).

**Variance of estimator** First, let's analyze the variance of our estimator. Note that our expectation is of the general form

$$\mu \triangleq \mathbb{E} [f(X)]$$

and that we estimate this quantity by generating  $n$  independent samples  $X_1, \dots, X_n \sim P(X)$  and computing

$$\hat{\mu} \triangleq \frac{1}{n} \sum_{i=1}^n f(X_i).$$

This is an unbiased estimator for  $\mu$  with error

$$[(\mu - \hat{\mu})^2] = \text{Var} [\hat{\mu}] = \frac{\text{Var} [\hat{\mu}]}{n},$$

which means that the error

$$\mu - \hat{\mu} = O\left(\sqrt{\frac{\text{Var} [\hat{\mu}]}{n}}\right)$$

and reducing it linearly requires a quadratic number of samples.

In our particular case, the function  $f$  is

$$f_{X=x}(Y) \triangleq L(X, Y) \nabla_{\lambda} \log q_{\lambda}(Y|X=x)$$

where we have made explicit that  $y$  is the random variable, and  $x$  is given.

**Control variates** Consider a function  $g(X)$  with known expectation

$$\mu_g \triangleq \mathbb{E}[g(X)]$$

Then we can define a new function  $\hat{f}$  such that

$$\hat{f}(X) \triangleq f(X) - g(X) + \mu_g.$$

This function is also an estimator for  $\mu$ , since

$$\begin{aligned} \mathbb{E}[\hat{f}(X)] &= \mathbb{E}[f(X)] - \mu_g + \mu_g \\ &= \mathbb{E}[f(X)], \end{aligned}$$

and a computation shows that the variance of the new function is

$$\begin{aligned} \text{Var}[\hat{f}(X)] &= \mathbb{E}[(f(X) - g(X) + \mu_g) - \mu]^2 \\ &= \mathbb{E}[(f(X) - g(X) + \mu_g)^2] - 2\mathbb{E}[(f(X) - g(X) + \mu_g)\mu] + \mathbb{E}[\mu^2] \\ &= \mathbb{E}[(f(X) - g(X) + \mu_g)^2] - 2\mathbb{E}[(f(X) - g(X) + \mu_g)]\mu + \mu^2 \\ &= \mathbb{E}[(f(X) - g(X) + \mu_g)^2] - 2\mu^2 + \mu^2 \\ &= \mathbb{E}[f(X)^2 + g(X)^2 + \mu_g^2 - 2f(X)g(X) + 2f(X)\mu_g - 2g(X)\mu_g] - \mu^2 \\ &= \mathbb{E}[f(X)^2] - \mathbb{E}[f(X)]^2 \\ &\quad - 2(\mathbb{E}[f(X)g(X)] - \mathbb{E}[f(X)]\mathbb{E}[g(X)]) \\ &\quad + \mathbb{E}[g(X)^2] - \mathbb{E}[g(X)]^2 \\ &= \text{Var}[f(X)] - 2\text{Cov}[f(X), g(X)] + \text{Var}[g(X)] \end{aligned}$$

This means we can get a reduction in variance whenever

$$\text{Cov}[f(X), g(X)] > \frac{1}{2} \text{Var}[g(X)].$$

The function  $g$  is called a *control variate*—it allows us to control the variance of  $f$ .

From the equality above we can see that this will be the case whenever  $f(X)$  and  $g(X)$  are strongly correlated. Our choice of control

variate will be made with the that in mind. Furthermore,  $\mathbb{E}[g(X)]$  must be known. What is an optimal control variate? Typically a control variate of the form  $ag$  is chosen with fixed, and  $a$  is optimized to maximize the correlation. This brings us to the generic formulation of a control variate:

$$\hat{f}(X) \triangleq f(X) - a(g(X) - \mathbb{E}[g(X)])$$

with variance

$$\text{Var}[\hat{f}(X)] = \text{Var}[f(X)] - 2a \text{Cov}[f(X), g(X)] + a^2 \text{Var}[g(X)]$$

We take a derivative of this with respect to  $a$

$$\frac{d}{da} \text{Var}[\hat{f}(X)] = -2 \text{Cov}[f(X), g(X)] + 2a \text{Var}[g(X)]$$

Setting this to zero and solving for  $a$  we obtain the optimal choice for  $a$

$$a = \frac{\text{Cov}[f(X), g(X)]}{\text{Var}[g(X)]}. \quad (\text{D.3})$$

Plugging in this solution into the expression for  $\text{Var}[\hat{f}(X)]$  and dividing by  $\text{Var}[f(X)]$  we get

$$\frac{\text{Var}[\hat{f}(X)]}{\text{Var}[f(X)]} = 1 - \frac{\text{Cov}[f(X), g(X)]}{\text{Var}[f(X)] \text{Var}[g(X)]} \quad (\text{D.4})$$

$$= 1 - \text{corr}^2[f(X), g(X)], \quad (\text{D.5})$$

which shows that given this choice of  $a$  the reduction in variance is directly determined by the correlation between  $f(X)$  and  $g(X)$ .

Bringing this all together, we let our new estimator be

$$\mathbb{E}[f(X)] = \mathbb{E}[\hat{f}(X)] \approx \frac{1}{n} \sum_{i=1}^n [f(X_i) - ag(X_i)] - \mu_g$$



**Example** (Ross, 2006) Suppose we want to use simulation to determine

$$\mathbb{E}[f(X)] = \mathbb{E}[e^X] = \int_0^1 e^x dx = e - 1$$

with  $X \sim \mathcal{U}(0, 1)$ . A natural control variate to use in this case is the random variable  $X$  itself:  $g(X) \triangleq X$ . We thus define the new estimator

$$\begin{aligned}\hat{f}(X) &= f(X) - g(X) + \mathbb{E}[g(X)] \\ &= e^X - X + \frac{1}{2}.\end{aligned}$$

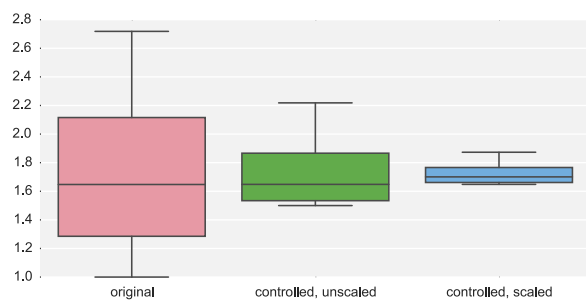
To compute the decrease in variance with this new estimator, we first note that

$$\begin{aligned}\text{Cov}(e^X, X) &= \mathbb{E}[Xe^X] - \mathbb{E}[X] \mathbb{E}[e^X] \\ &= \int_0^1 xe^x dx - \frac{e-1}{2} \\ &= 1 - \frac{e-1}{2} \approx 0.14086 \\ \text{Var}[e^X] &= \mathbb{E}[e^{2X}] - (\mathbb{E}[e^X])^2 \\ &= \int_0^1 e^{2x} dx - (1-e)^2 \\ &= \frac{e^2-1}{2} - (1-e)^2 \approx 0.2420 \\ \text{Var}[X] &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \\ &= \int_0^1 x^2 dx - \frac{1}{4} \\ &= \frac{1}{3} - \frac{1}{4} = \frac{1}{12}.\end{aligned}$$

When we choose  $a$  as in formula D.3 we can use formula D.4 to compute that

$$\begin{aligned}\frac{\text{Var}[\hat{f}(X)]}{\text{Var}[f(X)]} &= 1 - \frac{(0.14086)^2}{\frac{0.2420}{12}} \\ &\approx 0.0161.\end{aligned}$$

This is a reduction of 98.4 percent! A simulation illustrates what this looks like in practice with ... samples:



# E

---

## Syneval dataset

---

The following list is the complete set of categories that are evaluated in the Syneval dataset, taken from Marvin and Linzen (2018) with slight alterations. There is a slight ambiguity in the category of negative polarity items, with two different types on constructions, one of sentences starting with *most* and the other of sentences starting with *most*, combined into one positive set; we choose to separate them into two different categories. For the full list of lexical items used to build variants of these constructions we refer the reader to Marvin and Linzen (2018).

1. **Simple agreement:**
  - (a) The farmer *smiles*.
  - (b) \*The farmer *smile*.
2. **Agreement in a sentential complement:**
  - (a) The mechanics said the author *laughs*.
  - (b) \*The mechanics said the author *laugh*.
3. **Agreement in short VP coordination:**

- (a) The authors laugh and *swim*.
- (b) \*The authors laugh and *swims*.

**4. Agreement in long VP coordination:**

- (a) The author knows many different foreign languages and *enjoys* playing tennis with colleagues.
- (b) \*The author knows many different foreign languages and *enjoy* playing tennis with colleagues.

**5. Agreement across a prepositional phrase:**

- (a) The author next to the guards *smiles*.
- (b) \*The author next to the guards *smile*.

**6. Agreement across a subject relative clause:**

- (a) The author that likes the security guards *laughs*.
- (b) \*The author that likes the security guards *laugh*.

**7. Agreement across an object relative:**

- (a) The movies that the guard likes *are* good.
- (b) \*The movies that the guard likes *is* good.

**8. Agreement across an object relative (no *that*):**

- (a) The movies the guard likes *are* good.
- (b) \*The movies the guard likes *is* good.

**9. Agreement in an object relative:**

- (a) The movies that the guard *likes* are good.
- (b) \*The movies that the guard *like* are good.

**10. Agreement in an object relative (no *that*):**

(a) The movies the guard *likes* are good.

(b) \*The movies the guard *like* are good.

**11. Simple reflexive anaphora:**

(a) The author injured *himself*.

(b) \*The author injured *themselves*.

**12. Reflexive in sentential complement:**

(a) The mechanics said the author hurt *himself*.

(b) \*The mechanics said the author hurt *themselves*.

**13. Reflexive across a relative clause:**

(a) The author that the guards like injured *himself*.

(b) \*The author that the guards like injured *themselves*.

**14. Simple NPI:**

(a) *No* authors have ever been famous.

(b) \**Most* authors have ever been famous.

**15. Simple NPI (*the*):**

(a) *No* authors have ever been popular.

(b) \**The* authors have ever been popular.

**16. NPI across a relative clause:**

(a) *No* authors that *the* guards like have ever been famous.

(b) \**Most* authors that *no* guards like have ever been famous.

**17. NPI across a relative clause (*the*):**

(a) *No* authors that *the* guards like have ever been famous.

(b) \**The* authors that *no* guards like have ever been famous.

## References

---

- Andor, D., C. Alberti, D. Weiss, A. Severyn, A. Presta, K. Ganchev, S. Petrov, and M. Collins (2016). “Globally Normalized Transition-Based Neural Networks”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics. 2442–2452.
- Baker, J. K. (1979). “Trainable grammars for speech recognition”. *The Journal of the Acoustical Society of America*. 65(S1): S132–S132.
- Ballesteros, M., Y. Goldberg, C. Dyer, and N. A. Smith (2016). “Training with Exploration Improves a Greedy Stack LSTM Parser”. In: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Austin, Texas: Association for Computational Linguistics. 2005–2010.
- Bangalore, S. and A. K. Joshi (1999). “Supertagging: An approach to almost parsing”. *Computational linguistics*. 25(2): 237–265.
- Baydin, A. G., B. A. Pearlmutter, A. A. Radul, and J. M. Siskind (2018). “Automatic differentiation in machine learning: a survey”. *Journal of Machine Learning Research*. 18: 1–43.
- Bengio, Y., R. Ducharme, P. Vincent, and C. Jauvin (2003). “A neural probabilistic language model”. *Journal of machine learning research*. 3(Feb): 1137–1155.

- Black, E., S. Abney, D. Flickenger, C. Gdaniec, R. Grishman, P. Harrison, D. Hindle, R. Ingria, F. Jelinek, J. Klavans, *et al.* (1991). “A procedure for quantitatively comparing the syntactic coverage of English grammars”. In: *Speech and Natural Language: Proceedings of a Workshop Held at Pacific Grove, California, February 19-22, 1991*.
- Blei, D. M., A. Kucukelbir, and J. D. McAuliffe (2016). “Variational Inference: A Review for Statisticians”. *ArXiv e-prints*. Jan.
- Brennan, J. R., E. P. Stabler, S. E. Van Wagonen, W.-M. Luh, and J. T. Hale (2016). “Abstract linguistic structure correlates with temporal activity during naturalistic comprehension”. *Brain and language*. 157: 81–94.
- Buys, J. and P. Blunsom (2015a). “A Bayesian Model for Generative Transition-based Dependency Parsing”. In: *Proceedings of the Third International Conference on Dependency Linguistics, DepLing 2015, August 24-26 2015, Uppsala University, Uppsala, Sweden*. 58–67.
- Buys, J. and P. Blunsom (2015b). “Generative incremental dependency parsing with neural networks”. In: *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. Vol. 2. 863–869.
- Buys, J. and P. Blunsom (2018). “Neural Syntactic Generative Models with Exact Marginalization”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Vol. 1. 942–952.
- Carnie, A. (2010). *Constituent Structure*. *Oxford linguistics*. Oxford University Press. ISBN: 9780199583461.
- Caruana, R. (1997). “Multitask learning”. *Machine learning*. 28(1): 41–75.
- Chelba, C. and F. Jelinek (2000). “Structured language modeling”. *Computer Speech & Language*. 14(4): 283–332.
- Chelba, C., T. Mikolov, M. Schuster, Q. Ge, T. Brants, P. Koehn, and T. Robinson (2013). “One billion word benchmark for measuring progress in statistical language modeling”. *arXiv preprint arXiv:1312.3005*.

- Chelba, C., M. Norouzi, and S. Bengio (2017). “N-gram language modeling using recurrent neural network estimation”. *arXiv preprint arXiv:1703.10724*.
- Chen, S. F. and J. Goodman (1999). “An empirical study of smoothing techniques for language modeling”. *Computer Speech & Language*. 13(4): 359–394.
- Cheng, J., A. Lopez, and M. Lapata (2017). “A Generative Parser with a Discriminative Recognition Algorithm”. In: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Vancouver, Canada: Association for Computational Linguistics. 118–124.
- Chomsky, N. (1980). “Rules and representations”. *Behavioral and brain sciences*. 3(1): 1–15.
- Christiansen, M. H., C. M. Conway, and L. Onnis (2012). “Similar neural correlates for language and sequential learning: evidence from event-related brain potentials”. *Language and cognitive processes*. 27(2): 231–256.
- Collins, M. (2003). “Head-driven statistical models for natural language parsing”. *Computational linguistics*. 29(4): 589–637.
- Collobert, R. and J. Weston (2008). “A unified architecture for natural language processing: Deep neural networks with multitask learning”. In: *Proceedings of the 25th international conference on Machine learning*. ACM. 160–167.
- Collobert, R., J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa (2011). “Natural language processing (almost) from scratch”. *Journal of Machine Learning Research*. 12(Aug): 2493–2537.
- Conway, C. M. and D. B. Pisoni (2008). “Neurocognitive basis of implicit learning of sequential structure and its relation to language processing”. *Annals of the New York Academy of Sciences*. 1145(1): 113–131.
- Cross, J. and L. Huang (2016). “Span-Based Constituency Parsing with a Structure-Label System and Provably Optimal Dynamic Oracles”. In: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Austin, Texas: Association for Computational Linguistics. 1–11.



- Dozat, T. and C. D. Manning (2016). “Deep biaffine attention for neural dependency parsing”. *arXiv preprint arXiv:1611.01734*.
- Durrett, G. and D. Klein (2015). “Neural CRF Parsing”. In: *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)*. Beijing, China: Association for Computational Linguistics. 302–312.
- Dyer, C., A. Kuncoro, M. Ballesteros, and N. A. Smith (2016). “Recurrent Neural Network Grammars”. In: *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. San Diego, California: Association for Computational Linguistics. 199–209.
- Emami, A. and F. Jelinek (2005). “A neural syntactic language model”. *Machine learning*. 60(1-3): 195–227.
- Enguehard, É., Y. Goldberg, and T. Linzen (2017). “Exploring the Syntactic Abilities of RNNs with Multi-task Learning”. In: *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*. Vancouver, Canada: Association for Computational Linguistics. 3–14.
- Everaert, M. B., M. A. Huybregts, N. Chomsky, R. C. Berwick, and J. J. Bolhuis (2015). “Structures, not strings: linguistics as part of the cognitive sciences”. *Trends in cognitive sciences*. 19(12): 729–743.
- Finkel, J. R., A. Kleeman, and C. D. Manning (2008). “Efficient, Feature-based, Conditional Random Field Parsing”. In: *Proceedings of ACL-08: HLT*. Columbus, Ohio: Association for Computational Linguistics. 959–967.
- Frank, S. L., R. Bod, and M. H. Christiansen (2012). “How hierarchical is language use?” *Proceedings of the Royal Society B: Biological Sciences*. 279(1747): 4522–4531.
- Fried, D. and D. Klein (2018). “Policy Gradient as a Proxy for Dynamic Oracles in Constituency Parsing”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Melbourne, Australia: Association for Computational Linguistics. 469–476.

- Fu, M. C. (2006). “Gradient Estimation”. In: *Handbooks in Operations Research and Management Science (Volume 13)*. Ed. by B. L. N. Edited by Shane G. Henderson. Elsevier. Chap. 19. 575–616.
- Gaddy, D., M. Stern, and D. Klein (2018). “What’s Going On in Neural Constituency Parsers? An Analysis”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. New Orleans, Louisiana: Association for Computational Linguistics. 999–1010.
- Gallo, G., G. Longo, S. Pallottino, and S. Nguyen (1993). “Directed hypergraphs and applications”. *Discrete applied mathematics*. 42(2-3): 177–201.
- Giannakidou, A. (2011). “Negative and positive polarity items: Variation, licensing, and compositionality”. *The Handbook of Natural Language Meaning (second edition)*. Mouton de Gruyter, Berlin.
- Gillespie, M. and N. J. Pearlmutter (2011). “Hierarchy and scope of planning in subject–verb agreement production”. *Cognition*. 118(3): 377–397.
- Gillespie, M. and N. J. Pearlmutter (2013). “Against structural constraints in subject–verb agreement production.” *Journal of Experimental Psychology: Learning, Memory, and Cognition*. 39(2): 515.
- Glorot, X. and Y. Bengio (2010). “Understanding the difficulty of training deep feedforward neural networks”. In: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 249–256.
- Goldberg, Y. and J. Nivre (2013). “Training Deterministic Parsers with Non-Deterministic Oracles”. *Transactions of the Association for Computational Linguistics*. 1(Oct): 403–414.
- Goodman, J. (1999). “Semiring parsing”. *Computational Linguistics*. 25(4): 573–605.
- Gulordava, K., P. Bojanowski, E. Grave, T. Linzen, and M. Baroni (2018). “Colorless Green Recurrent Networks Dream Hierarchically”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. New Orleans, Louisiana: Association for Computational Linguistics. 1195–1205.

- Hale, J. (2001). “A probabilistic Earley parser as a psycholinguistic model”. In: *Proceedings of the second meeting of the North American Chapter of the Association for Computational Linguistics on Language technologies*. Association for Computational Linguistics. 1–8.
- Hale, J., C. Dyer, A. Kuncoro, and J. Brennan (2018). “Finding syntax in human encephalography with beam search”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics. 2727–2736.
- Hall, D., G. Durrett, and D. Klein (2014). “Less grammar, more features”. In: *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vol. 1. 228–237.
- He, H., J. Eisner, and H. Daume (2012). “Imitation learning by coaching”. In: *Advances in Neural Information Processing Systems*. 3149–3157.
- Hochreiter, S. and J. Schmidhuber (1997). “Long short-term memory”. *Neural computation*. 9(8): 1735–1780.
- Hockenmaier, J. and M. Steedman (2007). “CCGbank: a corpus of CCG derivations and dependency structures extracted from the Penn Treebank”. *Computational Linguistics*. 33(3): 355–396.
- Huddleston, R. D. and G. K. Pullum (2002). *The Cambridge Grammar of the English Language*. Cambridge University Press.
- Hudson, R. (2010). *An introduction to word grammar*. Cambridge University Press.
- Jang, E., S. Gu, and B. Poole (2017). “Categorical reparameterization with gumbel-softmax”. *ICLR*.
- Jelinek, F. (1997). “Information Extraction From Speech And Text, chapter 8”.
- Jordan, M., Z. Ghahramani, T. Jaakkola, and L. Saul (1999). “An Introduction to Variational Methods for Graphical Models”. *Machine Learning*. 37(2): 183–233.
- Jozefowicz, R., O. Vinyals, M. Schuster, N. Shazeer, and Y. Wu (2016). “Exploring the limits of language modeling”. *arXiv preprint arXiv:1602.02410*.

- Kalchbrenner, N., E. Grefenstette, and P. Blunsom (2014). “A convolutional neural network for modelling sentences”. *arXiv preprint arXiv:1404.2188*.
- Kingma, D. P. and J. Ba (2014). “Adam: A Method for Stochastic Optimization”. *CoRR*. abs/1412.6980.
- Kingma, D. P. and M. Welling (2014). “Auto-Encoding Variational Bayes”. In: *International Conference on Learning Representations*.
- Kitaev, N. and D. Klein (2018). “Constituency Parsing with a Self-Attentive Encoder”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics. 2676–2686.
- Klein, D. and C. D. Manning (2003). “Accurate Unlexicalized Parsing”. In: *Proceedings of the 41st Annual Meeting on Association for Computational Linguistics - Volume 1. ACL '03*. Sapporo, Japan: Association for Computational Linguistics. 423–430.
- Klein, D. and C. D. Manning (2004). “Parsing and hypergraphs”. In: *New developments in parsing technology*. Springer. 351–372.
- Kneser, R. and H. Ney (1995). “Improved backing-off for m-gram language modeling”. In: *icassp*. Vol. 1. 181e4.
- Kucukelbir, A., D. Tran, R. Ranganath, A. Gelman, and D. M. Blei (2017). “Automatic differentiation variational inference”. *The Journal of Machine Learning Research*. 18(1): 430–474.
- Kuncoro, A., M. Ballesteros, L. Kong, C. Dyer, G. Neubig, and N. A. Smith (2017). “What Do Recurrent Neural Network Grammars Learn About Syntax?” In: *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*. Valencia, Spain: Association for Computational Linguistics. 1249–1258.
- Kuncoro, A., C. Dyer, J. Hale, D. Yogatama, S. Clark, and P. Blunsom (2018). “LSTMs Can Learn Syntax-Sensitive Dependencies Well, But Modeling Structure Makes Them Better”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics. 1426–1436.

- Lafferty, J. D., A. McCallum, and F. C. N. Pereira (2001). “Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data”. In: *Proceedings of the Eighteenth International Conference on Machine Learning. ICML '01*. 282–289.
- Levy, R. (2008). “Expectation-based syntactic comprehension”. *Cognition*. 106(3): 1126–1177.
- Li, Z. and J. Eisner (2009). “First- and Second-Order Expectation Semirings with Applications to Minimum-Risk Training on Translation Forests”. In: *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*. Singapore: Association for Computational Linguistics. 40–51.
- Linzen, T., E. Dupoux, and Y. Goldberg (2016). “Assessing the Ability of LSTMs to Learn Syntax-Sensitive Dependencies”. *Transactions of the Association for Computational Linguistics*. 4: 521–535.
- Maddison, C. J., A. Mnih, and Y. W. Teh (2017). “The concrete distribution: A continuous relaxation of discrete random variables”. *ICLR*.
- Marcus, M. P., M. A. Marcinkiewicz, and B. Santorini (1993). “Building a large annotated corpus of English: The Penn Treebank”. *Computational linguistics*. 19(2): 313–330.
- Marcus, M., G. Kim, M. A. Marcinkiewicz, R. MacIntyre, A. Bies, M. Ferguson, K. Katz, and B. Schasberger (1994). “The Penn Treebank: annotating predicate argument structure”. In: *Proceedings of the workshop on Human Language Technology*. Association for Computational Linguistics. 114–119.
- Marvin, R. and T. Linzen (2018). “Targeted Syntactic Evaluation of Language Models”. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Brussels, Belgium: Association for Computational Linguistics. 1192–1202.
- McCoy, R. T., T. Linzen, and R. Frank (2018). “Revisiting the poverty of the stimulus: hierarchical generalization without a hierarchical bias in recurrent neural networks”. *arXiv preprint arXiv:1802.09091*.
- Merity, S., C. Xiong, J. Bradbury, and R. Socher (2016). “Pointer sentinel mixture models”. *arXiv preprint arXiv:1609.07843*.

- Miao, Y. and P. Blunsom (2016). “Language as a Latent Variable: Discrete Generative Models for Sentence Compression”. In: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Austin, Texas: Association for Computational Linguistics. 319–328.
- Mikolov, T., M. Karafiát, L. Burget, J. Černock, and S. Khudanpur (2010). “Recurrent neural network based language model”. In: *Eleventh Annual Conference of the International Speech Communication Association*.
- Mnih, A. and K. Gregor (2014). “Neural Variational Inference and Learning in Belief Networks”. In: *ICML*.
- Neubig, G., Y. Goldberg, and C. Dyer (2017a). “On-the-fly operation batching in dynamic computation graphs”. In: *Advances in Neural Information Processing Systems*. 3971–3981.
- Neubig, G., C. Dyer, Y. Goldberg, A. Matthews, W. Ammar, A. Anastopoulos, M. Ballesteros, D. Chiang, D. Clothiaux, T. Cohn, *et al.* (2017b). “Dynet: The dynamic neural network toolkit”. *arXiv preprint arXiv:1701.03980*.
- Nivre, J. (2005). “Dependency grammar and dependency parsing”. *MSI report*. 5133(1959): 1–32.
- Paisley, J., D. Blei, and M. Jordan (2012). “Variational Bayesian Inference with Stochastic Search”. In: *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*. Ed. by J. Langford and J. Pineau. *ICML ’12*. New York, NY, USA: Omnipress. 1367–1374.
- Pauls, A. and D. Klein (2012). “Large-scale syntactic language modeling with treelets”. In: *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*. Association for Computational Linguistics. 959–968.
- Peters, M. E., M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer (2018). “Deep contextualized word representations”. In: *Proc. of NAACL*.

- Petrov, S., L. Barrett, R. Thibaux, and D. Klein (2006). “Learning accurate, compact, and interpretable tree annotation”. In: *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*. Association for Computational Linguistics. 433–440.
- Ranganath, R., S. Gerrish, and D. Blei (2014). “Black box variational inference”. In: *Artificial Intelligence and Statistics*. 814–822.
- Rennie, S. J., E. Marcheret, Y. Mroueh, J. Ross, and V. Goel (2017). “Self-Critical Sequence Training for Image Captioning”. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*: 1179–1195.
- Rezende, D. J., S. Mohamed, and D. Wierstra (2014). “Stochastic Backpropagation and Approximate Inference in Deep Generative Models”. In: *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*. 1278–1286.
- Roark, B. (2001). “Probabilistic top-down parsing and language modeling”. *Computational linguistics*. 27(2): 249–276.
- Rosenfeld, R. (1996). “A maximum entropy approach to adaptive statistical language modelling”. *Computer Speech Language*. 10(3): 187–228. ISSN: 0885-2308.
- Ross, S. M. (2006). *Simulation, Fourth Edition*. Orlando, FL, USA: Academic Press, Inc.
- Schulman, J., N. Heess, T. Weber, and P. Abbeel (2015). “Gradient estimation using stochastic computation graphs”. In: *Advances in Neural Information Processing Systems*. 3528–3536.
- Sekine, S. and M. Collins (1997). “EVALB bracket scoring program”. URL: <http://www.cs.nyu.edu/cs/projects/proteus/evalb>.
- Sennrich, R. (2017). “How Grammatical is Character-level Neural Machine Translation? Assessing MT Quality with Contrastive Translation Pairs”. In: *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*. Valencia, Spain: Association for Computational Linguistics. 376–382.
- Smith, N. A. (2012). “Adversarial evaluation for models of natural language”. *arXiv preprint arXiv:1207.0245*.

- Søgaard, A. and Y. Goldberg (2016). “Deep multi-task learning with low level tasks supervised at lower layers”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Berlin, Germany: Association for Computational Linguistics. 231–235.
- Stern, M., J. Andreas, and D. Klein (2017a). “A Minimal Span-Based Neural Constituency Parser”. In: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vancouver, Canada: Association for Computational Linguistics. 818–827.
- Stern, M., D. Fried, and D. Klein (2017b). “Effective Inference for Generative Neural Parsing”. In: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics. 1695–1700.
- Swayamdipta, S., S. Thomson, K. Lee, L. Zettlemoyer, C. Dyer, and N. A. Smith (2018). “Syntactic Scaffolds for Semantic Structures”. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Brussels, Belgium: Association for Computational Linguistics. 3772–3782.
- Tesnière, L. (1959). “Elements of structural syntax, translated by Timothy Osborne and Sylvain Kahane”.
- Titov, I. and J. Henderson (2007). “A Latent Variable Model for Generative Dependency Parsing”. In: *Proceedings of the 10th International Conference on Parsing Technologies. IWPT '07*. Prague, Czech Republic: Association for Computational Linguistics. 144–155. ISBN: 978-1-932432-90-9.
- Tran, K., A. Bisazza, and C. Monz (2018). “The Importance of Being Recurrent for Modeling Hierarchical Structure”. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Brussels, Belgium: Association for Computational Linguistics. 4731–4736.
- Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin (2017). “Attention is all you need”. In: *Advances in Neural Information Processing Systems*. 5998–6008.



- Vlachos, A. (2013). “An investigation of imitation learning algorithms for structured prediction”. In: *European Workshop on Reinforcement Learning*. 143–154.
- Wainwright, M. J. and M. I. Jordan (2008). “Graphical Models, Exponential Families, and Variational Inference”. *Found. Trends Mach. Learn.* 1(1-2): 1–305.
- Warstadt, A., A. Singh, and S. R. Bowman (2018). “Neural Network Acceptability Judgments”. *arXiv preprint arXiv:1805.12471*.
- Williams, R. J. (1992). “Simple statistical gradient-following algorithms for connectionist reinforcement learning”. *Machine learning*. 8(3-4): 229–256.
- Wilson, A. C., R. Roelofs, M. Stern, N. Srebro, and B. Recht (2017). “The marginal value of adaptive gradient methods in machine learning”. In: *Advances in Neural Information Processing Systems*. 4148–4158.
- Xiang, M., B. Dillon, and C. Phillips (2009). “Illusory licensing effects across dependency types: ERP evidence”. *Brain and Language*. 108(1): 40–55.
- Yin, P., C. Zhou, J. He, and G. Neubig (2018). “StructVAE: Tree-structured Latent Variable Models for Semi-supervised Semantic Parsing”. *arXiv preprint arXiv:1806.07832*.
- Zaremba, W., I. Sutskever, and O. Vinyals (2014). “Recurrent neural network regularization”. *arXiv preprint arXiv:1409.2329*.
- Zhang, Y. and D. Weiss (2016). “Stack-propagation: Improved Representation Learning for Syntax”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics. 1557–1566.
- Zweig, G. and C. J. Burges (2011). “The microsoft research sentence completion challenge”. *Tech. rep.* Citeseer.