

Natural Language Processing

Cristian Perez Jensen

August 3, 2024

Note that these are not the official lecture notes of the course, but only notes written by a student of the course. As such, there might be mistakes. The source code can be found at github.com/cristianpjensen/eth-cs-notes. If you find a mistake, please create an issue or open a pull request.

Contents

1	Backpropagation	1
2	Log-linear modeling	4
2.1	Softmax	6
2.2	The exponential family	7
3	Multi-layer perceptrons	8
3.1	Word embeddings	8
3.2	Sentiment analysis	9
4	Structured prediction	10
5	Language modeling	11
5.1	n-grams	12
5.2	Recurrent neural networks	12
6	Semirings	14
7	Part-of-speech tagging	16
7.1	Conditional random fields	16
8	Finite-state automata	18
8.1	WFST composition	18
8.2	Pathsum	19
8.3	Lehmann's algorithm	20
9	Transliteration	22
10	Constituency parsing	23
10.1	Context-free grammars	23
10.2	Parsing	24
10.3	Cocke-Kasami-Younger algorithm	25
11	Dependency parsing	27
11.1	Chu-Liu-Edmonds algorithm	28
12	Semantic parsing	30
12.1	Linear-indexed grammars	30
12.2	Lambda calculus	31
12.3	Combinatory logic	32
12.4	Combinatory categorial grammars	33
13	Transformers	35
13.1	Translation	35

List of symbols

\doteq	Equality by definition
\approx	Approximate equality
\propto	Proportional to
\mathbb{N}	Set of natural numbers
\mathbb{R}	Set of real numbers
$i : j$	Set of natural numbers between i and j . <i>I.e.</i> , $\{i, i+1, \dots, j\}$
$f : A \rightarrow B$	Function f that maps elements of set A to elements of set B
$\mathbb{1}\{\text{predicate}\}$	Indicator function (1 if predicate is true, otherwise 0)
$\mathbf{v} \in \mathbb{R}^n$	n -dimensional vector
$\mathbf{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
$\mathbf{T} \in \mathbb{R}^{d_1 \times \dots \times d_n}$	Tensor
\mathbf{M}^\top	Transpose of matrix \mathbf{M}
\mathbf{M}^{-1}	Inverse of matrix \mathbf{M}
$\det(\mathbf{M})$	Determinant of \mathbf{M}
$\frac{d}{dx}f(x)$	Ordinary derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}$
$\frac{\partial}{\partial x}f(\mathbf{x})$	Partial derivative of $f(\mathbf{x})$ w.r.t. x at point $\mathbf{x} \in \mathbb{R}^n$
$\nabla_{\mathbf{x}}f(\mathbf{x}) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\mathbf{x} \in \mathbb{R}^n$
$\nabla_{\mathbf{x}}^2f(\mathbf{x}) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\mathbf{x} \in \mathbb{R}^n$
$\boldsymbol{\theta} \in \Theta$	Parametrization of a model, where Θ is a compact subset of \mathbb{R}^K
\mathcal{X}	Input space
\mathcal{Y}	Output space
$\mathcal{D} \subseteq \mathcal{X} \times \mathcal{Y}$	Labeled training data

1 Backpropagation

Backpropagation is the single most important algorithm in modern machine learning, because it is used to compute gradients of composite functions efficiently. It is a linear-time dynamic program for computing derivatives, *i.e.*, it stores intermediate results to be as fast as forward propagation.

In machine learning, most of the time, we have inputs $x \in \mathcal{X}$ and outputs $y \in \mathcal{Y}$ from a dataset $\mathcal{D} \subseteq \mathcal{X} \times \mathcal{Y}$, and we want to fit some function $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ such that it minimizes some loss function,

$$\sum_{(x,y) \in \mathcal{D}} \ell(f_\theta(x), y).$$

We need this function's gradient to be able to use an algorithm such as gradient descent to optimize it.¹ For a composite function, it is time consuming to derive the gradient by hand. Thus, we use backpropagation to automatically compute the gradients, as long as we have access to the derivatives of its primitive functions.

¹ Most of the time, this function cannot be solved in closed form.

Example 1 (Computation graph). A composite function can be represented using intermediate variables such that each variable is computed by a single primitive function. Let's say we have the following function,

$$f(x, y) = \sin(xy + \exp(y)).$$

Then, we can represent the intermediate variables as the following,

$$\begin{aligned} z_1 &= xy \\ z_2 &= \exp(y) \\ z_3 &= z_1 + z_2 \\ z_4 &= \sin(z_3). \end{aligned}$$



We could also describe a function as a labeled, directed acyclic² hypergraph³, where each node is a variable and each hyperedge is labeled with a function. Given a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $y = f(x) \in \mathbb{R}^m$, for input

² The fact that our computation graph is acyclic makes it possible for backpropagation to be linear.

³ A hypergraph allows the edges to have multiple sources and targets. This is needed because functions can have multiple inputs and multiple outputs.

x_j and outputs y_i , Bauer's formula gives the following,

$$\frac{\partial y_i}{\partial x_j} = \sum_{p \in P(j,i)} \prod_{(k,l) \in p} \frac{\partial z_l}{\partial z_k},$$

where $P(j,i)$ is the set of paths from vertex j to vertex i and $p \in P(j,i)$ is the set of edges that make up the path p . *I.e.*, the partial derivative is a sum over all paths in the computation graph, where the derivative over each path is computed by the chain rule.⁴ When computing the partial derivatives for functions with dense computation graphs naively, we are typically summing over an exponential number of paths, because many of these partial derivatives are recomputed many times. We can use dynamic programming to store these values and avoid computing them again. In this case, the amount of computation scales linearly with the number of edges.

$$^4 \frac{d}{dx} f(g(x)) = \frac{d}{dg(x)} f(g(x)) \cdot \frac{d}{dx} g(x).$$

```

1: function FORWARDPROPAGATION( $f, x$ )
2:    $z_i \leftarrow \begin{cases} x_i & \text{if } i \leq m \\ 0 & \text{otherwise} \end{cases}$   $\triangleright$  Initialize input variables
3:   for  $i = m + 1, \dots, n$  do
4:      $z_i \leftarrow g_i(z_{\text{Parents}(i)})$   $\triangleright$  Set intermediate variables
5:   end for
6:   return  $z$ 
7: end function

```

Algorithm 1. Forward propagation algorithm that assumes that the edges are topologically sorted so $i < j$ implies that z_i is computed before z_j .

```

1: function BACKPROPAGATION( $f, x$ )
2:    $z \leftarrow \text{FORWARDPROPAGATION}(f, x)$ 
3:    $\frac{\partial f}{\partial z_i} \leftarrow \begin{cases} 1 & \text{if } i = n \\ 0 & \text{otherwise} \end{cases}$   $\triangleright$  Base case
4:   for  $i = n - 1, \dots, 1$  do  $\triangleright O(n)$ 
5:      $\frac{\partial f}{\partial z_i} \leftarrow \sum_{z_p \in \text{Parents}(z_i)} \frac{\partial f}{\partial z_p} \frac{\partial z_p}{\partial z_i}$   $\triangleright$  Chain rule
6:   end for
7:   return  $\nabla_x f$ 
8: end function

```

Algorithm 2. Backpropagation algorithm that assumes that the edges are topologically sorted so $i < j$ implies that z_i is computed before z_j .

The general framework that any backpropagation framework uses is the following,

1. Write down a composite function as a hypergraph with intermediate variables as nodes and hyperedges labeled with primitive functions;
2. Given a set of inputs, perform forward propagation through the graph to compute the function's value (Algorithm 1);
3. Run backpropagation on the graph using the stored forward values (Algorithm 2). Intuitively, we set up a dynamic programming table using Bauer's formula.

Example 2 (Backpropagation table for Example 1).

$$\frac{\partial f}{\partial z_4} = 1$$

$$\frac{\partial f}{\partial z_3} = \frac{\partial f}{\partial z_4} \frac{\partial z_4}{\partial z_3} = \cos(z_3)$$

$$\frac{\partial f}{\partial z_2} = \frac{\partial f}{\partial z_3} \frac{\partial z_3}{\partial z_2} = \cos(z_3)$$

$$\frac{\partial f}{\partial z_1} = \frac{\partial f}{\partial z_3} \frac{\partial z_3}{\partial z_1} = \cos(z_3)$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial z_1} \frac{\partial z_1}{\partial x} = \cos(z_3)y$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial z_1} \frac{\partial z_1}{\partial y} + \frac{\partial f}{\partial z_2} \frac{\partial z_2}{\partial y} = \cos(z_3)x + \cos(z_3)\exp(y).$$

2 Log-linear modeling

Let's say we want to model the conditional probability $p(y | x)$. A naive way of doing this is the following,

$$p(y | x) \doteq \frac{\text{count}(x, y)}{\text{count}(x)}.$$

There are two main problems with this interpretation of discrete conditional probability,

- Suppose $\text{count}(x, y) = 0$, then the probability will be 0, *i.e.*, the model says that y is impossible in context x ;
- There is no way to look at finer-grained aspects of x , *i.e.*, some values of x might be related.

Thus, we need a more general framework for modeling conditional distributions. One such general framework is to simply exponentiate some scoring function $\text{score} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ that we construct,⁵ and let the conditional probability be proportional to it,

$$p(y | x) \propto \exp \text{score}(x, y) > 0.$$

The linear scoring function looks like the following,

$$\text{score}(x, y) = \theta^\top f(x, y)$$

with feature weights $\theta \in \mathbb{R}^K$ and $f(x, y) \in \mathbb{R}^K$ as a vector describing y in context x . The conditional probability then looks like the following,

$$p_\theta(y | x) = \frac{1}{Z_\theta} \exp(\theta^\top f(x, y))$$

$$Z_\theta(x) = \sum_{y' \in \mathcal{Y}} \exp(\theta^\top f(x, y')).$$

This is called log-linear modeling, because if we take the logarithm of the conditional probability, we get a linear model,

$$\log p_\theta(y | x) = \theta^\top f(x, y) - \log Z_\theta(x).$$

The design of the *feature function* $f(x, y)$ is a big portion of the work in log-linear modeling. It can be split into two parts: preprocessing and extracting features. The preprocessing simply consists of steps such as tokenization, lower-casing, stemming, stop-word removal, and reducing vocabulary. After the preprocessing, we can obtain features. Examples include one-hot encoding, bag-of-words, n -grams, and word embeddings.

Maximum likelihood estimation (MLE) is a way of finding the parameters $\theta \in \Theta$ that minimizes the *negative log-likelihood* (a.k.a. *cross entropy*) of the training data, *i.e.*, we want to minimize the following,

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} - \log \prod_{(x, y) \in \mathcal{D}} p_\theta(y | x)$$

$$= \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{(x, y) \in \mathcal{D}} -\log p_\theta(y | x),$$

⁵ The exponentiation makes sure that it is non-negative.

where Θ is a compact (bounded and closed) subset of \mathbb{R}^K .⁶ In log-linear modeling, this objective function is convex, thus any local minimum is a global minimum. We usually optimize the log-likelihood with gradient-based methods,

⁶ The reason for not taking $\theta \in \mathbb{R}^K$ is that the weights will likely go to infinity.

$$\begin{aligned}
\nabla_{\theta} \sum_{(x,y) \in \mathcal{D}} -\log p_{\theta}(y | x) &= \sum_{(x,y) \in \mathcal{D}} \nabla_{\theta} \left(\log Z_{\theta}(x) - \theta^{\top} f(x, y) \right) \\
&= \sum_{(x,y) \in \mathcal{D}} \nabla_{\theta} \log Z_{\theta}(x) - \nabla_{\theta} \theta^{\top} f(x, y) \\
&= \sum_{(x,y) \in \mathcal{D}} \frac{1}{Z_{\theta}(x)} \sum_{y' \in \mathcal{Y}} \nabla_{\theta} \exp \left(\theta^{\top} f(x, y') \right) - f(x, y) \\
&= \sum_{(x,y) \in \mathcal{D}} \sum_{y' \in \mathcal{Y}} \frac{1}{Z_{\theta}(x)} \exp \left(\theta^{\top} f(x, y') \right) \nabla_{\theta} \theta^{\top} f(x, y') - f(x, y) \\
&= \sum_{(x,y) \in \mathcal{D}} \sum_{y' \in \mathcal{Y}} p_{\theta}(y' | x) f(x, y') - f(x, y) \\
&= \sum_{(x,y) \in \mathcal{D}} \mathbb{E}_{y'}[f(x, y')] - \sum_{(x,y) \in \mathcal{D}} f(x, y).
\end{aligned}$$

Due to convexity, the global minimum is the only point that has its gradient equal 0. Thus, at the optimal parameters, the following is the case,

$$f(x, y) = \mathbb{E}_{y'}[f(x, y')].$$

Therefore, the optimum is where the observed feature counts $f(x, y)$ look like the expected feature counts $\mathbb{E}_{y'}[f(x, y')]$ through the lens of the model. In other words, the training data looks exactly like what our model predicts through the eyes of our feature function. This is referred to as *expectation matching*.

Furthermore, we can derive the Hessian of the negative log-likelihood to be the covariance matrix of $f(x, \cdot)$ w.r.t. θ ,

$$\begin{aligned}
\nabla_{\theta^\top} p_\theta(y | x) &= \nabla_{\theta^\top} Z_\theta(x)^{-1} \exp(\theta^\top f(x, y)) \\
&= Z_\theta(x)^{-1} \nabla_{\theta^\top} \exp(\theta^\top f(x, y)) + \left(\nabla_{\theta^\top} Z_\theta(x)^{-1} \right) \exp(\theta^\top f(x, y)) \\
&= \frac{1}{Z_\theta(x)} \exp(\theta^\top f(x, y)) \nabla_{\theta^\top} \theta^\top f(x, y) - \frac{1}{Z_\theta(x)^2} \exp(\theta^\top f(x, y)) \nabla_{\theta^\top} Z_\theta(x) \\
&= p_\theta(y | x) f(x, y)^\top - p_\theta(y | x) \frac{1}{Z_\theta(x)} \sum_{y' \in \mathcal{Y}} \nabla_{\theta^\top} \exp(\theta^\top f(x, y')) \\
&= p_\theta(y | x) f(x, y)^\top - p_\theta(y | x) \frac{1}{Z_\theta(x)} \sum_{y' \in \mathcal{Y}} \exp(\theta^\top f(x, y')) \nabla_{\theta^\top} \theta^\top f(x, y') \\
&= p_\theta(y | x) f(x, y)^\top - p_\theta(y | x) \sum_{y' \in \mathcal{Y}} p_\theta(y' | x) f(x, y')^\top \\
&= p_\theta(y | x) \left(f(x, y)^\top - \mathbb{E}_{y'} [f(x, y')^\top] \right).
\end{aligned}$$

$$\begin{aligned}
\nabla_\theta^2 - \log p_\theta(y | x) &= \sum_{(x, y) \in \mathcal{D}} \nabla_{\theta^\top} \nabla_\theta - \log p_\theta(y | x) \\
&= \sum_{(x, y) \in \mathcal{D}} \nabla_{\theta^\top} \sum_{y' \in \mathcal{Y}} p_\theta(y' | x) f(x, y') - f(x, y) \\
&= \sum_{(x, y) \in \mathcal{D}} \sum_{y' \in \mathcal{Y}} f(x, y') \nabla_{\theta^\top} p_\theta(y' | x) \\
&= \sum_{(x, y) \in \mathcal{D}} \sum_{y' \in \mathcal{Y}} f(x, y') p_\theta(y' | x) \left(f(x, y')^\top - \mathbb{E}_{y''} [f(x, y'')^\top] \right) \\
&= \sum_{(x, y) \in \mathcal{D}} \sum_{y' \in \mathcal{Y}} p_\theta(y' | x) f(x, y') f(x, y')^\top - \mathbb{E}_{y''} [f(x, y'')^\top] \sum_{y' \in \mathcal{Y}} p_\theta(y' | x) f(x, y') \\
&= \sum_{(x, y) \in \mathcal{D}} \mathbb{E}_{y'} [f(x, y) f(x, y)^\top] - \mathbb{E}_{y'} [f(x, y)^\top] \mathbb{E}_{y'} [f(x, y)] \\
&= \sum_{(x, y) \in \mathcal{D}} \text{Cov}_{y' \sim p_\theta(\cdot | x)} [f(x, y')].
\end{aligned}$$

2.1 Softmax

The softmax : $\mathbb{R}^K \rightarrow \Delta^{K-1}$ function is the default way of building probabilistic models using neural networks, because it maps vectors to categorical probability distributions. It is basically the same as log-linear modeling. It is defined as

$$\text{softmax}(\mathbf{h})_y \doteq \frac{\exp(h_y/T)}{\sum_{y' \in \mathcal{Y}} \exp(h_{y'}/T)}.$$

Usually, the temperature is set to $T = 1$.⁷ This is a generalization of log-linear modeling, where, instead of $\theta^\top f(x, y)$, we can use any function of the input.

The probability simplex Δ^{K-1} is a subspace of $\mathbb{R}_{\geq 0}^K$ such that the sum of the components of its elements is 1. It is denoted as Δ^{K-1} , because it has $K - 1$ free parameters, and looks like a triangle in three dimensions.

⁷ As $T \rightarrow 0$, softmax becomes the argmax function and as $T \rightarrow \infty$, softmax becomes a uniform categorical distribution.

2.2 The exponential family

The *exponential family* is a family of probability distributions of the following form,

$$p_{\theta}(x) = \frac{1}{Z_{\theta}} h(x) \exp(\theta^{\top} \phi(x)),$$

where Z_{θ} is the partition function that normalizes the probability distribution, $h(x)$ determines the support of the function, θ are the canonical parameters, and $\phi(x)$ are the sufficient statistics. Importantly, $h(x)$ may not depend on θ . Any distribution that can be brought into this form is part of this family.⁸

The advantage of the exponential family is that they have conjugate priors, which make intractable posteriors tractable. Intuitively, this is because the posterior must be in the same form as the prior, thus we must be able to summarize the data into a finite vector (sufficient statistics).

⁸ Notice that the form of the exponential family is a generalization of softmax.

3 Multi-layer perceptrons

For log-linear models to find an appropriate model, the data has to be linearly separable.⁹ The solution to this problem is the *multi-layer perceptron* (MLP) [Haykin, 1994]. They jointly learn a non-linear feature function with the model's parameters. MLPs consist of n alternating linear projections and non-linearities,¹⁰

$$\begin{aligned} h_n &= \sigma_n(\mathbf{W}_n^\top h_{n-1}) \\ h_1 &= \sigma_1(\mathbf{W}_1^\top e(x)), \end{aligned}$$

where $e(x) \in \mathbb{R}^d$ is a continuous vector encoding of the input x . Then, we can combine this non-linear feature representation of the input with the parameters to obtain a categorical probability distribution,

$$\text{softmax}(\boldsymbol{\theta}^\top h_n).$$

Basically, the only difference is that we now learn the feature function.

Using backpropagation, we can compute the derivative of a weight matrix \mathbf{W}_k of an MLP as follows,

$$\begin{aligned} \frac{\partial \ell}{\partial \mathbf{W}_k} &= \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial h_n} \left(\prod_{m=k+1}^n \frac{\partial h_m}{\partial h_{m-1}} \right) \frac{\partial h_k}{\partial \mathbf{W}_k} \\ &= \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial h_n} \left(\prod_{m=k+1}^n \frac{\partial}{\partial h_{m-1}} \sigma_m(\mathbf{W}_m^\top h_{m-1}) \right) \frac{\partial}{\partial \mathbf{W}_k} \sigma_k(\mathbf{W}_k^\top h_{k-1}) \\ &= \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial h_n} \left(\prod_{m=k+1}^n \sigma'_m(\mathbf{W}_m^\top h_{m-1}) \mathbf{W}_m \right) \sigma'_k(\mathbf{W}_k^\top h_{k-1}) h_{k-1} \end{aligned}$$

If we use the ReLU activation function, this becomes the following,

$$= \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial h_n} \left(\prod_{m=k+1}^n \mathbb{1}\{\mathbf{W}_m^\top h_{m-1} > 0\} \cdot \mathbf{W}_m \right) \mathbb{1}\{\mathbf{W}_k^\top h_{k-1} > 0\} \cdot h_{k-1}.$$

This is the cause of “dead” neurons. When their values become negative, their derivative becomes 0. Thus, there is no way for gradient descent to learn anymore. Also, the derivative of all children of “dead” neurons become 0, because of the chain rule.

3.1 Word embeddings

To be able to use MLPs, we would need a continuous vector encoding of words/sentences. A naive idea would be to encode one-hot word counts of sentences. However, this approach discards word ordering information. A naive solution would be to then encode n -grams to regain some word ordering information, however the vectors would explode in size to $\mathcal{O}(|V|^n)$, where $|V|$ is the vocabulary size.

A great idea is to use unsupervised learning to train word embeddings on large corpora.¹¹ The first model that made use of this idea is Skip-

⁹ If the data is non-linear, we can use a feature function f that makes the feature vector linearly separable by hacking the non-linearity into the feature function (e.g., if the data is ellipsoidal, we would choose an ellipsoidal feature function.) But, this requires us to know the decision boundary's shape a priori to set the feature function.

¹⁰ The non-linearities are very important, because otherwise we would just get a linear model, since a stack of linear transformations are equal to some single linear transformation. Thus, if the model would not contain non-linearities, we would not gain any expressiveness.

¹¹ Easy to get due to the internet.

Gram [Mikolov et al., 2013]. Its high-level idea is that it predicts whether a word w' is in the context of a word w .¹² The weights of this model are the word embeddings and they are used to predict the model's objective. The idea is then that the model needs a good representation of the words to be able to do this task successfully, thus we can use the embeddings that this model trains for other tasks.

The first step of Skip-Gram is to preprocess the corpus. This is done by collecting positive and negative samples.¹³ To collect positive samples, it goes through all words w and collects all words w' in the context of w . Then, it randomly generates other words w' as negative samples, *i.e.*, they are not in the context of w , and adds them to the dataset.

Then, we use the embeddings to predict whether a context word is in a word's context,

$$p(c | w) = \frac{1}{Z(w)} \exp(e_w(w)^\top e_c(c))$$

$$Z(w) = \sum_{c \in V} \exp(e_w(w)^\top e_c(c)).$$

We use two different embedding types, because if $w = c$, then $e(w)^\top e(c)$ would be positive and thus the probability very high, even though this case is very unlikely, because words do usually not appear in their own context. Thus, to mitigate this, we use two different context and word embeddings.

3.2 Sentiment analysis

An example application is sentiment analysis [Pang et al., 2008], which classifies sentences as positive or negative. MLPs and word embeddings can be used for this and work very well [Iyyer et al., 2015]. The model first embeds all the words in the sentence and pools them together into one vector representation of the sentence.¹⁴ Then, this vector is passed into an MLP, which is used as input to a softmax with a single output that tells us how positive the sentence is.

¹² "You shall know a word by the company it keeps."

¹³ A reason for collecting negative samples is because then we do not need to compute the normalizing constant $Z(w)$. This is due to the fact that we are simply maximizing the output for positive samples and minimizing the output for negative samples while training. But, if we did not have negative samples, we would have to normalize the output to be a probability between 0 and 1, because otherwise we would not know whether the output is good or not.



Figure 1. Architecture of a simple sentiment classifier, where f_θ is a multi-layer perceptron.

¹⁴ Pooling together discards word order, but for this simple task, it will still work very well.

4 Structured prediction

In NLP, we often have the case that the input and output of a model has some structure. *E.g.*, in context-free parsing, we have a sentence as input and want to output a parse tree. To be able to train a model and perform inference, we need to be able to compute the log-likelihood,

$$p_{\theta}(y \mid x) \doteq \frac{1}{Z_{\theta}(x)} \exp(\text{score}_{\theta}(y, x))$$

$$Z_{\theta}(x) \doteq \sum_{y' \in \mathcal{Y}} \exp(\text{score}_{\theta}(y', x)).$$

However, often there are an exponential, or even infinite, amount of possible structures $y \in \mathcal{Y}$ for an input $x \in \mathcal{X}$. For training, this has the result that the normalizer Z_{θ} is very inefficient to compute, since it is a sum over a very large amount of values. For inference, this has the result that we would need to search a very large space for the best output y .

The solution to this problem is that we need to make use of the structure to design algorithms for computing the normalizer and finding the most probable output. The next sections will include NLP problems, where we need to design algorithms to compute the normalizer efficiently.

5 Language modeling

In language modeling, Σ is a finite, non-empty set of symbols. In the context of natural language, this is usually set to be the vocabulary of words or tokens. A string over an alphabet Σ is any finite sequence of alphabet symbols. The output space \mathcal{Y} is set to the set of all possible strings Σ^* , which is infinite, thus there is no way to sum over every possible structure, nor is there an easy way to output the maximum scoring structure.

A *language model* (LM) is a probability distribution over Σ^* , i.e., LMs assign probabilities to strings $\mathbf{y} \in \Sigma^*$. We can discriminate between two types of LMs,

- *Globally normalized* LMs, which define a single scoring function $\text{score}_\theta : \Sigma^* \rightarrow \mathbb{R}$ and normalizes the scores across all $\mathbf{y} \in \Sigma^*$,¹⁵

$$p(\mathbf{y}) \doteq \frac{1}{Z_\theta} \exp \text{score}_\theta(\mathbf{y});$$

- *Locally normalized* LMs, which decompose string probabilities into conditional probabilities $p(y_i \mid \mathbf{y}_{<i})$ over symbols y_i given the previous symbols $\mathbf{y}_{<i}$,

$$p(\mathbf{y}) \doteq p(\text{EOS} \mid \mathbf{y}) \cdot \prod_{i=1}^N p(y_i \mid \mathbf{y}_{<i})$$

$$p(y_i \mid \mathbf{y}_{<i}) \doteq \frac{1}{Z_\theta(\mathbf{y}_{<i})} \exp \text{score}_\theta(y_i, \mathbf{y}_{<i}).$$

Local LMs are collections of conditional probability distributions $p(y_i \mid \mathbf{y}_{<i})$, which intuitively tells us how probable symbol y_i is to follow the already seen string $\mathbf{y}_{<i}$. In practice, we also need beginning-of-sentence (BOS) and end-of-sentence (EOS) symbols. We condition on BOS for the first token to model that the sequence starts with y_1 and we condition on the entire string with EOS to model the probability that no more symbols follow.

A well-defined LM always assigns probability to EOS given any history $\mathbf{y}_{<i}$, because otherwise we could get in the situation where we have a history that never ends in EOS, which will not result in a string. Models that have this problem are called non-tight. The probability of all sentences in a non-tight model do not sum to 1. To mitigate this problem, we ensure a model is tight by forcing $p(\text{EOS} \mid \mathbf{y}_{<i}) > \xi > 0$ for every history $\mathbf{y}_{<i}$.

The problem with local LMs is that there are infinitely many distributions $p(\cdot \mid \mathbf{y}_{<i})$ with $\mathbf{y}_{<i} \in \Sigma^*$. Thus, we need some way of being able to compute all these distributions.

¹⁵ Global LMs are not used much, because their normalizer requires a sum over an infinite set. Thus, we will focus on local LMs.

We can also use local LMs to generate strings by continuously sampling from the probability distribution $p(y_i \mid \mathbf{y}_{<i})$ until we sample EOS, e.g.,

$$\text{HE WALKS THE} \begin{cases} p(\text{DOG} \mid \text{HE} \dots) = 0.5 \\ p(\text{CAT} \mid \text{HE} \dots) = 0.24 \\ p(\text{EOS} \mid \text{HE} \dots) = 0.01. \end{cases}$$

5.1 n -grams

The n -gram solution is to limit histories $\mathbf{y}_{<i}$ to a length $n - 1$. This assumption leads to the following probability distribution,

$$p(y_i | \mathbf{y}_{<i}) = p(y_i | y_{i-n+1}, \dots, y_{i-1}).$$

This ensures a finite number $|\Sigma|^{n-1}$ of histories.

The naive implementation of n -gram would be to define a separate conditional probability distribution for every possible context of size $n - 1$. Thus, we can assign each history a probability distribution based on the counts we observe in training data,

$$p(y_i | y_{i-n+1}, \dots, y_{i-1}) = \frac{\text{count}(y_{i-n+1}, \dots, y_{i-1}, y_i)}{\text{count}(y_{i-n+1}, \dots, y_{i-1})}.$$

However, this does not allow us to share parameters between histories, which might be very similar and give much insight.¹⁶ Furthermore, it is very memory inefficient, since we need to store $(|\bar{\Sigma}| - 1) |\Sigma|^{n-1}$ parameters to specify the $|\Sigma|^{n-1}$ conditional distributions for each history, where $\bar{\Sigma} = \Sigma \cup \{\text{EOS}\}$.

The neural n -gram [Bengio et al., 2000] is a more efficient approach that does allow parameter sharing. It uses word embeddings to encode the words and histories,

$$p(y_i | \mathbf{y}_{<i}) = \frac{\exp(\mathbf{e}(y_i)^\top \mathbf{h}_i)}{\sum_{y' \in \Sigma} \exp(\mathbf{e}(y')^\top \mathbf{h}_i)},$$

where $\mathbf{h}_i = \text{enc}(\mathbf{y}_{<i}) = \text{enc}(y_{i-n+1}, \dots, y_{i-1})$.¹⁷ The crucial idea of this approach is that if the word embedding is similar to the encoding of the history, it is more likely. This approach has $\mathcal{O}(d|V|)$ space complexity for all n .

5.2 Recurrent neural networks

The fixed history size of the n -gram model is not realistic. We want to be able to encode the entire history, which is possible with *recurrent neural networks* (RNN). RNNs work by having two inputs at every timestep: the time-dependent hidden state \mathbf{h}_{i-1} , representing all words before the current one, and embedding of the current token $\mathbf{e}(y_{i-1})$. These are combined to represent the entire history. There are many variations, but the simplest one is the Elman RNN [Elman, 1990],

$$\mathbf{h}_i = \sigma(\mathbf{W}_h \mathbf{h}_{i-1} + \mathbf{W}_x \mathbf{e}(y_{i-1}) + \mathbf{b}),$$

where $\mathbf{W}_h \in \mathbb{R}^{d \times d}$, $\mathbf{W}_x \in \mathbb{R}^{d \times d}$, and $\mathbf{b} \in \mathbb{R}^d$ are learned parameters. At their core, RNNs are just a non-linear combination of the recurrent state and the inputs.

¹⁶ E.g., the distributions over

$$p(\cdot | \text{SHE WALKS}) \quad p(\cdot | \text{HE WALKS})$$

are parametrized independently, even though their distributions should be extremely similar.

¹⁷ Bengio et al. [2000] used a neural network to encode the history. In this architecture, the words in the history $\mathbf{y}_{i-n+1:i-1}$ are concatenated (preserve word-order) and used as input to a neural network that outputs a d -dimensional representation \mathbf{h}_i .

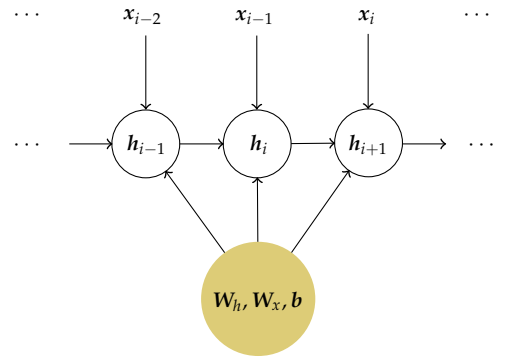


Figure 2. Diagram of the RNN architecture. Each hidden state \mathbf{h}_i has “seen” all previous tokens $x_{1:i-1}$.

RNNs are trained by unrolling the timesteps and applying backpropagation. The problem with this is that RNNs become prone to the vanishing/exploding gradient problem, because computing the gradient of the parameters involves a lot of matrix operations with itself. The LSTM [Hochreiter and Schmidhuber, 1997] and GRU [Cho et al., 2014] architectures seek to solve this problem.

$$\begin{aligned}
\frac{\partial \ell}{\partial \mathbf{W}_h} &= \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial \mathbf{h}_N} \frac{\partial \mathbf{h}_N}{\partial \mathbf{W}_h} \\
&= \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial \mathbf{h}_N} \sum_{i=1}^N \frac{\partial \mathbf{h}_N}{\partial \mathbf{h}_i} \frac{\partial \mathbf{h}_i}{\partial \mathbf{W}_h} \\
&= \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial \mathbf{h}_N} \sum_{i=1}^N \left(\prod_{m=i+1}^N \frac{\partial \mathbf{h}_m}{\partial \mathbf{h}_{m-1}} \right) \frac{\partial \mathbf{h}_i}{\partial \mathbf{W}_h} \\
&= \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial \mathbf{h}_N} \sum_{i=1}^N \left(\prod_{m=i+1}^N \frac{\partial}{\partial \mathbf{h}_{m-1}} \sigma(\mathbf{W}_h \mathbf{h}_{i-1} + \mathbf{W}_x \mathbf{e}(y_{i-1}) + \mathbf{b}) \right) \frac{\partial}{\partial \mathbf{W}_h} \sigma(\mathbf{W}_h \mathbf{h}_{i-1} + \mathbf{W}_x \mathbf{e}(y_{i-1}) + \mathbf{b}) \\
&= \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial \mathbf{h}_N} \sum_{i=1}^N \left(\prod_{m=i+1}^N \sigma'(\mathbf{W}_h \mathbf{h}_{i-1} + \mathbf{W}_x \mathbf{e}(y_{i-1}) + \mathbf{b}) \mathbf{W}_h \right) \sigma'(\mathbf{W}_h \mathbf{h}_{i-1} + \mathbf{W}_x \mathbf{e}(y_{i-1}) + \mathbf{b}) \mathbf{h}_{i-1}
\end{aligned}$$

If we use the tanh activation function with a derivative that never exceeds 1, we get the following,

$$\approx \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial \mathbf{h}_N} \sum_{i=1}^N \mathbf{h}_{i-1} \prod_{m=i+1}^N \mathbf{W}_h.$$

This contains many multiplications of \mathbf{W}_h with itself, which causes exploding or vanishing gradient (dependent on whether the determinant of \mathbf{W}_h is greater than or less than 1).

6 Semirings

Definition 3 (Monoid). A monoid is a 3-tuple $\langle \mathbb{K}, \odot, e \rangle$, such that

1. \odot is associative for all values in \mathbb{K} . $\forall x, y, z \in \mathbb{K}$,

$$(x \odot y) \odot z = x \odot (y \odot z);$$

2. $e \in \mathbb{K}$ is the identity element. $\forall x \in \mathbb{K}$,

$$x \odot e = x.$$

Definition 4 (Semiring). A semiring is a 5-tuple $\langle \mathbb{K}, \oplus, \otimes, 0, 1 \rangle$, such that

1. $\langle \mathbb{K}, \oplus, 0 \rangle$ is a commutative monoid;
2. $\langle \mathbb{K}, \otimes, 1 \rangle$ is a monoid;
3. \otimes distributes over \oplus . $\forall x, y, z \in \mathbb{K}$,

$$\begin{aligned} (x \oplus y) \otimes z &= (x \otimes z) \oplus (y \otimes z) \\ z \otimes (x \oplus y) &= (z \otimes x) \oplus (z \otimes y); \end{aligned}$$

4. 0 is an annihilator for \otimes . $\forall x \in \mathbb{K}$,

$$\begin{aligned} 0 \otimes x &= 0 \\ x \otimes 0 &= 0. \end{aligned}$$

Semirings are very useful for generalizing algorithms that only make use of associativity, commutativity, and distributivity. For example, if we have an efficient algorithm for computing the normalizer,

$$Z = \sum_{y \in \mathcal{Y}} \prod_{n=1}^N \exp \text{score}(y_n).$$

We can “semiringify” it to compute the following,

$$\bigoplus_{y \in \mathcal{Y}} \bigotimes_{n=1}^N \exp \text{score}(y_n).$$

Then, we can use any semiring with this algorithm. For inference, we would then want to use for example the Viterbi semiring $\langle \mathbb{R}, \max, \times, 0, 1 \rangle$ to compute the following,

$$\max_{y \in \mathcal{Y}} \prod_{n=1}^N \exp \text{score}(y_n).$$

Definition 5 (Closed semiring). A closed semiring is a semiring with an additional unary operation: the Kleene star,

$$x^* = \bigoplus_{n=0}^{\infty} x^{\otimes n}.$$

The Kleene star must obey the following two axioms,

$$x^* = \mathbf{1} \oplus x \otimes x^*$$

$$x^* = \mathbf{1} \oplus x^* \otimes x.$$

Closedness allows us to compute infinite sums within a semiring. For example, the real semiring is closed if we let its set be $(-1, 1)$, because then we can compute the Kleene star to be the following using the closed form of the geometric series,

$$x^* = \sum_{n=0}^{\infty} x^n = \frac{1}{1-x}.$$

7 Part-of-speech tagging

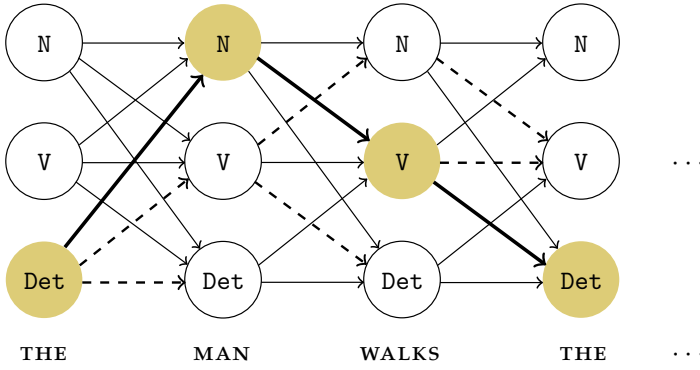


Figure 3. Example POS graph with $\mathcal{T} = \{N, V, \text{Det}\}$. For inference, we want to find the best N -length path within this graph. For training, we want to compute the sum over all N -length paths within this graph. The dashed edges are the backpointers to the argmaxes, and the thick line is the best path.

In *part-of-speech tagging* (POS tagging), we want to predict a POS tag $t \in \mathcal{T}$ for every word of the input sentence w of length N . In other words, given an N -dimensional input sequence of words $w \in \Sigma^N$, we want to output an N -dimensional sequence of tags $t \in \mathcal{T}^N$. This can be seen as searching through a POS graph as in Figure 3. The output space \mathcal{T}^N is exponential, so we need to design an algorithm to efficiently compute the normalizer $Z(w)$, and find the maximum scoring tagging.

7.1 Conditional random fields

Conditional random fields (CRF) are a conditional probabilistic model for structured labeling. Whereas a classifier predicts a label for a single sample without considering neighboring samples in the structure, a CRF does take context into account. In other words, CRFs are a model for computing the normalizer Z in a structured labeling case.

In the sequence labeling case of POS tagging, we will assume that tags only depend on their immediate neighbors,

$$\text{score}(t, w) = \sum_{n=1}^N \text{score}(\langle t_{n-1}, t_n \rangle, w, n),$$

which can be further decomposed into transition and emission scores,

$$\text{score}(\langle t_{n-1}, t_n \rangle, w, n) = \text{transition}(\langle t_{n-1}, t_n \rangle) + \text{emission}(w_n, t_n).$$

This balances how likely t_n is to follow t_{n-1} and how likely word w_n is to be assigned tag t_n .

We can use the new decomposed scoring function to find an efficient

Note that the bigram assumption does not mean that the current tag only depends on the previous and current word, because the word representations can be anything. *E.g.*, if we use a bidirectional RNN for the word representations, the tags will still depend on the entire input sentence. However, a problem that this can cause is that it cannot correctly tag garden-path sentences, since we cannot change the start of the tagging after seeing the end of the tagging. An example garden-path sentence is “The horse raced past the barn fell.”

algorithm for computing the normalizer under a semiring,

$$\begin{aligned}
& \bigoplus_{t \in \mathcal{T}^N} \bigotimes_{n=1}^N \exp \text{score}(\langle t_{n-1}, t_n \rangle, w, n) \\
&= \bigoplus_{t_1 \in \mathcal{T}} \cdots \bigoplus_{t_n \in \mathcal{T}} \exp \text{score}(\langle t_0, t_1 \rangle, w, 1) \otimes \cdots \otimes \exp \text{score}(\langle t_{N-1}, t_N \rangle, w, N) \\
&= \bigoplus_{t_1 \in \mathcal{T}} \exp \text{score}(\langle t_0, t_1 \rangle, w, 1) \otimes \left(\cdots \otimes \left(\bigoplus_{t_n \in \mathcal{T}} \exp \text{score}(\langle t_{N-1}, t_N \rangle, w, N) \right) \right),
\end{aligned}$$

where distributivity is used in the last step. The backward and forward algorithms (Algorithms 3 and 4) are direct results of this rederivation, which compute the normalizer in $\mathcal{O}(N \cdot |\mathcal{T}|^2)$.

```

1: function BACKWARDALGORITHM( $w, \text{score}, \langle A, \oplus, \otimes, 0, 1 \rangle$ )
2:    $\beta[N] \leftarrow 1$ 
3:   for  $n = N - 1, \dots, 0$  do
4:     for  $t_n \in \mathcal{T}$  do
5:        $\beta[n, t_n] \leftarrow \bigoplus_{t_{n+1} \in \mathcal{T}} \exp(\text{score}(\langle t_n, t_{n+1} \rangle), w, n + 1) \otimes$ 
         $\beta[n + 1, t_{n+1}]$ 
6:     end for
7:   end for
8:   return  $\beta[0, \text{BOT}]$ 
9: end function

```

Algorithm 3. Backward algorithm that computes the semiring-sum over all taggings of a sentence w . It can be seen as iteratively computing the larger semiring-sum from the derived equation.

```

1: function FORWARDALGORITHM( $w, \text{score}, \langle A, \oplus, \otimes, 0, 1 \rangle$ )
2:    $\alpha[0] \leftarrow 1$ 
3:   for  $n = 1, \dots, N$  do
4:     for  $t_n \in \mathcal{T}$  do
5:        $\alpha[n, t_n] \leftarrow \bigoplus_{t_{n-1} \in \mathcal{T}} \exp(\text{score}(\langle t_{n-1}, t_n \rangle), w, n) \otimes$ 
         $\alpha[n - 1, t_{n-1}]$ 
6:     end for
7:   end for
8:   return  $\alpha[N, \text{EOT}]$ 
9: end function

```

Algorithm 4. Forward algorithm that computes the same thing as Algorithm 3, but then in a different fashion. This version is more intuitive, because it starts at the beginning and ends at the end.

We can thus use these algorithms to compute the normalizer $Z(w)$ during training. For inference, the Viterbi algorithm is a version of the backward algorithm under the Viterbi semiring, where backpointers to the argmax are kept such that the maximally scoring tagging can be constructed by backtracking the backpointer.

8 Finite-state automata

A *finite-state automaton* (FSA) is a computational device that determines whether a string $s \in \Sigma^*$ is an element of a given language $L \subseteq \Sigma^*$. To check whether a string is part of a language defined by an FSA, the FSA reads in letters of an input string $s \in \Sigma^*$. Then, it transitions from state to state according to the transition function δ and the letters $a \in s$. If there is a path from an initial state to a final state while taking transitions as specified, the FSA accepts the string and is part of its language.

Definition 6 (Finite-state automaton). A finite-state automaton is a 5-tuple $\mathcal{A} = \langle \Sigma, Q, I, F, \delta \rangle$, such that

- Σ is an alphabet;
- Q is a finite set of states;
- $I \subseteq Q$ is the set of initial states;
- $F \subseteq Q$ is the set of final states;
- $\delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times Q$ is a finite multi-set that defines the transitions between states. Let $\langle q_0, a, q_1 \rangle \in \delta$ be a transition, then, if we make that transition, we go to state q_1 from q_0 and concatenate a to the current string.

A *weighted finite-state automaton* (WFSA) adds weights from a semiring to the transitions ($\delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times \mathbb{K} \times Q$), initial states ($\lambda : I \rightarrow \mathbb{K}$), and final states ($\rho : F \rightarrow \mathbb{K}$). Weights are added together using the \otimes operator.

A *weighted finite-state transducer* (WFST) further adds an output alphabet. The transitions then go from state to state while mapping input characters to output characters. Formally, $\delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times (\Omega \cup \{\epsilon\}) \times \mathbb{K} \times Q$, where Ω is the output alphabet.

8.1 WFST composition

WFST composition of two transducers \mathcal{T}_1 and \mathcal{T}_2 is a common operation that involves mapping the inputs of \mathcal{T}_1 to the outputs of \mathcal{T}_2 . This requires the output alphabet of \mathcal{T}_1 to be equal to the input alphabet of \mathcal{T}_2 . Intuitively, it is the same as first running the input through \mathcal{T}_1 then that output through \mathcal{T}_2 ,

$$x \xrightarrow{\mathcal{T}_1} z \xrightarrow{\mathcal{T}_2} y.$$

The weight of mapping x to y using the composition of two WFSTs is the semiring-sum of the weights of all possible transformations of the above form.

An n -gram LM can be represented by a WFSA by setting the states to be the history. The initial states are all `bos`. We can only end on the `eos` state. A state can go to another if the history can follow the other and the weight is the probability of that happening. The initial and final weights are all 1.

A CRF can also be represented by a WFSA, where each POS tag is a state. We can start and end on any state. The transitions then go from tag to tag, where the weight is the score of the target tag following the source tag. The initial weights are then the score of the tag following `bos`, and the target weights are the score of `eos` following the tag.

The Kleene star of an alphabet Σ is defined as

$$\Sigma^* \doteq \bigcup_{n=0}^{\infty} \Sigma^n.$$

Definition 7 (WFST composition). Formally, the composition $\mathcal{T} = \mathcal{T}_1 \circ \mathcal{T}_2$ of two WFSTs,

$$\begin{aligned}\mathcal{T}_1 &= \langle \Sigma, \Omega, Q_1, I_1, F_1, \delta_1, \lambda_1, \rho_1 \rangle \\ \mathcal{T}_2 &= \langle \Omega, \Xi, Q_2, I_2, F_2, \delta_2, \lambda_2, \rho_2 \rangle,\end{aligned}$$

is the WFST $\mathcal{T} = \langle \Sigma, \Xi, Q, I, F, \delta, \lambda, \rho \rangle$, such that

$$\mathcal{T}(x, y) = \bigoplus_{z \in \Omega^*} \mathcal{T}_1(x, z) \otimes \mathcal{T}_2(z, y).$$

```

1: function NAIVECOMPOSITION( $\mathcal{T}_1, \mathcal{T}_2$ )
2:    $\mathcal{T} \leftarrow \langle \Sigma, \Omega, Q, I, F, \delta, \lambda, \rho \rangle$       ▷ Create a new WFST
3:   for  $q_1, q_2 \in Q_1 \times Q_2$  do
4:     for  $q_1 \xrightarrow{a:b/w_1} q'_1, q_2 \xrightarrow{c:d/w_2} q'_2 \in E_1(q_1) \times E_2(q_2)$  do
5:       if  $b = c$  then
6:          $Q \leftarrow Q \cup \{(q_1, q_2), (q'_1, q'_2)\}$       ▷ Add states
7:          $\delta \leftarrow \delta \cup \{(q_1, q_2) \xrightarrow{a:d/w_1 \otimes w_2} (q'_1, q'_2)\}$       ▷ Add arcs
8:       end if
9:     end for
10:  end for
11:  for  $(q_1, q_2) \in Q$  do      ▷ Initial and final weights
12:     $\lambda((q_1, q_2)) \leftarrow \lambda_1(q_1) \otimes \lambda_2(q_2)$ 
13:     $\rho((q_1, q_2)) \leftarrow \rho_1(q_1) \otimes \rho_2(q_2)$ 
14:  end for
15:  return  $\mathcal{T}$ 
16: end function

```

Algorithm 5. Naive version of the algorithm for computing the composition of two WFSTs.

8.2 Pathsum

A path $\pi \in \delta^*$ is an ordered set of *consecutive* transitions,

$$\left(q_1 \xrightarrow{a_1:b_1/w_1} q_2, q_2 \xrightarrow{a_2:b_2/w_2} q_3, \dots, q_{N-1} \xrightarrow{a_N:b_N/w_N} q_N \right).$$

The weight of this path is defined as

$$w(\pi) \doteq \lambda(q_1) \otimes \bigotimes_{n=1}^N w_n \otimes \rho(q_N).$$

The pathsum is then the semiring-sum over all paths in a WFST,

$$\begin{aligned}Z(\mathcal{T}) &\doteq \bigoplus_{\pi \in \Pi(\mathcal{A})} w(\pi) \\ &= \bigoplus_{\pi \in \Pi(\mathcal{A})} \left(\lambda(q_1) \otimes \bigotimes_{n=1}^N w_n \otimes \rho(q_N) \right).\end{aligned}$$

Under the real semiring, the pathsum computes

$$\sum_{\pi \in \Pi(\mathcal{A})} \left(\lambda(q_1) \times \prod_{n=1}^N w_n \times \rho(q_N) \right),$$

which is the normalizer, while under the Viterbi semiring, the pathsum computes

$$\max_{\pi \in \Pi(\mathcal{A})} \left(\lambda(q_1) \times \prod_{n=1}^N w_n \times \rho(q_N) \right),$$

which is the maximum score of a path.

8.3 Lehmann's algorithm

Generally, there are an infinite amount of paths in a WFST, because of possible cycles. Thus, we need to design an algorithm that can compute this potentially infinite semiring-sum. Lehmann's algorithm [Lehmann, 1977] computes the semiring-sum matrix R over all inner paths between all pairs of nodes in $\mathcal{O}(|Q|^3)$ under a closed semiring. Then, we can use these to compute the normalizer with the following equation,

$$Z(\mathcal{T}) \doteq \bigoplus_{i,k \in Q} \lambda(q_i) \otimes R_{ik} \otimes \rho(q_k).$$

```

1: function LEHMANN( $W, \langle \mathcal{W}, \oplus, \otimes, \mathbf{0}, \mathbf{1} \rangle$ )           ▷  $\mathcal{W}$  is closed
2:    $R^{(0)} \leftarrow W$                                        ▷  $W \in \mathcal{W}^{|Q| \times |Q|}$ 
3:   for  $j \leftarrow 1, \dots, |Q|$  do
4:     for  $i \leftarrow 1, \dots, |Q|$  do
5:       for  $k \leftarrow 1, \dots, |Q|$  do
6:          $R_{ik}^{(j)} \leftarrow R_{ik}^{(j-1)} \oplus R_{ij}^{(j-1)} \otimes (R_{jj}^{(j-1)})^* \otimes R_{jk}^{(j-1)}$ 
7:       end for
8:     end for
9:   end for
10:  return  $I \oplus R^{(|Q|)}$ 
11: end function

```

Algorithm 6. Lehmann's algorithm to compute the inner path semiring-sums.

Lehmann's algorithm is a dynamic programming algorithm where $R_{ik}^{(j)}$ is the semiring-sum over all paths between i and k through $\{1, \dots, j\}$. The base case is then the weight matrix W since that contains the weights of going from i to k directly.

Intuitively, the recurrence relationship,

$$R_{ik}^{(j)} = R_{ik}^{(j-1)} \oplus R_{ij}^{(j-1)} \otimes (R_{jj}^{(j-1)})^* \otimes R_{jk}^{(j-1)},$$

says that the sum over all paths through $\{1, \dots, j\}$ is equal to the sum over all paths through $\{1, \dots, j-1\}$ plus the sum over all paths through j . The recurrence here is that we can use $R_{ik}^{(j-1)}$ as the sum over all paths

through $\{1, \dots, j-1\}$. We can compute the sum over all paths through j by the following,

$$R_{ij}^{(j-1)} \otimes \left(R_{jj}^{(j-1)}\right)^* \otimes R_{jk}^{(j-1)},$$

which is the weight of all paths from i to j , j to j (cycles), and finally j to k .

```

1: function FLOYD-WARSHALL( $W$ )
2:    $R^{(0)} \leftarrow W$ 
3:   for  $j \leftarrow 1, \dots, |Q|$  do
4:     for  $i \leftarrow 1, \dots, |Q|$  do
5:       for  $k \leftarrow 1, \dots, |Q|$  do
6:          $R_{ik}^{(j)} \leftarrow \min\{R_{ik}^{(j-1)}, R_{ij}^{(j-1)} + R_{jk}^{(j-1)}\}$ 
7:       end for
8:     end for
9:   end for
10:  return  $\min\{I, R^{(|Q|)}\}$ 
11: end function

```

Algorithm 7. Floyd-Warshall algorithm to compute the shortest path distance between any two vertices in a graph without negative cycles. This is very similar to Lehmann's algorithm in the semiring $\langle \mathbb{R}, \min, +, \infty, 0 \rangle$.

The Floyd-Warshall algorithm [Floyd, 1962, Warshall, 1962] has the same intuition as Lehmann's algorithm. It is also a dynamic program with the same state space. Also, the recurrence relationship is very similar, but it assumes no negative cycles, thus it does not need the Kleene star to account for cycles from j to j .

9 Transliteration

Transliteration is the mapping of strings in one character set to strings in another character set. An example of this is the phonetic translation of English words. Formally, we want to develop a probabilistic model that can map strings from input vocabulary Σ to an output vocabulary Ω , *i.e.*, we want to compute $p(y | x)$ for all $x \in \Sigma^*, y \in \Omega^*$. We can use a WFST to specify the transliteration of Σ^* to Ω^* as a globally normalized model. The scoring function is then the semiring-sum over all paths that aligns x with y ,

$$\text{score}(y, x) \doteq \sum_{\pi \in \Pi(y)} w(\pi).$$

To compute the normalizer, we need to design a WFST such that the input can only be x and that the output can be any element of Ω^* . To compute $\text{score}(y, x)$, we need a WFST such that the input can only be x and the output only y . We can guarantee such behavior by defining three transducers,

- \mathcal{T}_x is the transducer that maps x to x ;
- \mathcal{T}_θ is the transducer that maps any source string in Σ^* to any target string Ω^* (Figure 4);
- \mathcal{T}_y is the transducer that maps y to y .

We can compose $\mathcal{T}_x \circ \mathcal{T}_\theta$ to get a transducer that has as input only x and output any target string in Ω^* . We can use Lehmann's algorithm to compute the normalizer $Z_\theta(x)$ using the real semiring, and the maximally scoring output using the Viterbi semiring. We can then use the transducer composition $\mathcal{T}_x \circ \mathcal{T}_\theta \circ \mathcal{T}_y$ to compute $\text{score}_\theta(y, x)$ by using the real semiring (or the log semiring to keep it in log-space). Thus, we have all the components we need for training and inference.



Figure 4. Mapping WFST between the alphabets $\{A, B, C\}$ and $\{\alpha, \beta, \gamma\}$. The vertex label denotes the last added symbol to the output string. This WFST maps any string of the input alphabet to any string of the output alphabet. Intuitively, it has insertion, substitution, and deletion operations.

10 Constituency parsing



Figure 5. Two possible constituency trees of the ambiguous sentence "I shot an elephant in my pajamas."

A parse tree is a hierarchy of constituents, where a constituent is a multi-word unit that functions as a single unit. Each constituent encapsulates all of its leaf descendants, which are the words of the sentence. We say that a tree yields the sentence that can be found on its leaves. However, language is ambiguous, so some sentences have multiple trees that yield it. The goal is to compute the best parse tree that yields a given natural language input sentence.

10.1 Context-free grammars

Intuitively, a grammar defines a set of sentences that are deemed grammatical. Any sentence that can be yielded by a tree that consists of rules defined by the grammar is deemed grammatical.

Definition 8 (Context-free grammar). A context-free grammar is a 4-tuple $\langle \mathcal{N}, S, \Sigma, \mathcal{R} \rangle$, such that

- \mathcal{N} is a set of non-terminal symbols, written as uppercase letters N_1, N_2, \dots ;
- $S \in \mathcal{N}$ is a distinguished start non-terminal. Every complete parse tree must have this symbol at its root;
- Σ is an alphabet of terminal symbols, written as lowercase letters a_1, a_2, \dots ;
- \mathcal{R} is a set of production rules of the following form,

$$N \rightarrow \alpha,$$

where $N \in \mathcal{N}$ and $\alpha \in (\mathcal{N} \cup \Sigma)^*$.

A *context-free grammar* (CFG) encodes a subset of Σ^* , where a sentence is only part of the subset if we can construct a tree from \mathcal{R} that yields the sentence, starting from S .

10.2 Parsing

We might be able to assign multiple trees to a single sentence. To be able to pick the best tree, we can assign a probability to each rule, and pick the tree with the highest probability.

Definition 9 (Probabilistic context-free grammar). A probabilistic CFG is a 5-tuple $G = \langle \mathcal{N}, S, \Sigma, \mathcal{R}, p \rangle$, where $p : \mathcal{R} \rightarrow [0, 1]$ is a locally normalized probability distribution over rules. The probability of a tree under a PCFG is defined as follows,

$$p(t \mid s) = \prod_{r \in t} p(r).$$

Instead of probabilities, we could also, more generally, assign weights to each production rule. The score of assigning a tree t to a sentence w then decomposes over the production rules,

$$\text{score}(t, w) \doteq \sum_{r \in t} \text{score}(r),$$

where a tree t is simply a multiset of rules. However, we run into the problem that the normalizer $Z(w)$ will diverge if the ruleset contains a cycle rule, e.g., $N \rightarrow N$.


```

1: function WEIGHTEDCKY( $w, \langle \mathcal{N}, \mathcal{S}, \Sigma, \mathcal{R} \rangle, \text{score}$ )
2:    $\mathbf{C} \leftarrow \mathbf{0}$  ▷ Chart
3:   for  $i = 1, \dots, N$  do
4:     for  $X \rightarrow w_i \in \mathcal{R}$  do
5:        $\mathbf{C}[i, i + 1, X] \leftarrow \mathbf{C}[i, i + 1, X] \oplus \text{exp score}(X \rightarrow w_i)$ 
6:     end for
7:   end for
8:   for  $\ell = 2, \dots, N$  do
9:     for  $i = 1, \dots, N - \ell + 1$  do
10:       $k \leftarrow i + \ell$ 
11:      for  $j = i + 1, \dots, k - 1$  do
12:        for  $X \rightarrow YZ \in \mathcal{R}$  do
13:           $\mathbf{C}[i, k, X] \leftarrow \mathbf{C}[i, k, X] \oplus \text{exp score}(X \rightarrow YZ) \otimes$ 
             $\mathbf{C}[i, j, Y] \otimes \mathbf{C}[j, k, Z]$ 
14:        end for
15:      end for
16:    end for
17:  end for
18:  return  $\mathbf{C}[1, N + 1, \mathcal{S}]$ 
19: end function

```

Algorithm 8. Semiringified CKY algorithm that runs in $\mathcal{O}(N^3|\mathcal{R}|)$, where N is the input string length. If all possible rules exist, the runtime is $\mathcal{O}(N^3|\mathcal{N}|^3)$.

starts from the bottom of the chart and works its way up using dynamic programming.

This algorithm runs in $\mathcal{O}(N^3 \cdot |\mathcal{R}|)$, where N is the length of the sentence. We can semiringify this algorithm to compute the best parse, where the Viterbi semiring finds the maximally scoring tree, and the real semiring computes the normalizer.

11 Dependency parsing

Dependency parsing is an alternative to constituency parsing. The basic idea is to link every word with its *syntactic head*.¹⁸



¹⁸ We encode this as an arc in a graph, see Figures 7 and 8.

Figure 7. Projective dependency tree of “The boy eats Rösti.”

In a dependency tree, only one word gets to be the root, and each word has a single parent, called its *syntactic head*. This allows for words to be linked that have other words in between not part of this structure, in contrast to constituency parsings.

There are two types of dependency trees, projective and non-projective. Projective dependency trees do not allow for crossing arcs, which make them closely related to constituency trees. Non-projective dependency trees do allow for crossing arcs, which will be the focus of this text. An example of a non-projective dependency tree can be found in Figure 8.



Figure 8. Non-projective dependency tree, without labels, of “Ryan ate Rösti, which was delicious.”

As always, we want to be able to parametrize a probability distribution over non-projective spanning trees, given a sentence w . So, we need to be able to compute the normalizer $Z(w)$. However, there are $(N - 1)^{N-2}$ spanning trees with the single root constraint, thus we need to design a more efficient algorithm that makes use of its structure. We do this by decomposing the scoring function over the edges,

$$\text{score}(t, w) \doteq \text{score}(r, w) + \sum_{(i \rightarrow j) \in t} \text{score}(i, j, w).$$

Thus, we only need a matrix A , containing $\exp \text{score}(i, j, w)$ and a vector ρ , containing $\exp \text{score}(r, w)$.

Theorem 12 (Kirchhoff’s matrix-tree theorem [Kirchhoff, 1847]). For an undirected unweighted graph \mathcal{G} with N vertices, let L be the graph Laplacian,

$$L_{ij} = \begin{cases} -A_{ij} & i \neq j \\ \sum_{k \neq i} A_{kj} & \text{otherwise,} \end{cases}$$

where A is the adjacency matrix, *i.e.*, $A_{ij} = 1$ if $i \sim j$, otherwise $A_{ij} = 0$. Let $\hat{L}_i \in \mathbb{R}^{(N-1) \times (N-1)}$ be the matrix created by removing the i -th row and column of L . Then, we have

$$N_T(\mathcal{G}) = \det(\hat{L}_i),$$

where $N_T(\mathcal{G})$ is the number of trees in \mathcal{G} .

Tutte [1948] generalized Kirchhoff’s MTT to directed trees, which allows us to compute the normalizer $Z(w)$ in $\mathcal{O}(N^3)$ as follows,¹⁹

$$Z(w) = \det(L).$$

However, this does not account for the single-root constraint. Koo et al. [2007] further generalized Tutte’s MTT by modifying the graph Laplacian,

$$L_{ij} = \begin{cases} \rho_j & i = 1 \\ -A_{ij} & i \neq j \\ \sum_{k \neq i} A_{kj} & \text{otherwise.} \end{cases}$$

So, we can now compute $Z(w) = \det(L)$ in $\mathcal{O}(N^3)$.

However, we are not able to semiringify this algorithm. Thus, we must design a new algorithm for inference.

11.1 Chu-Liu-Edmonds algorithm

A valid dependency tree must adhere to the following three constraints,

- All non-root nodes have exactly one incoming edge;
- No cycles;
- Only one outgoing edge from the root.

This is called an *arborescence*,²⁰ and we can find the maximum-weight arborescence with the Chu-Liu-Edmonds algorithm [Chu, 1965, Edmonds et al., 1967], sped up to $\mathcal{O}(N^2)$ by Tarjan [1977].

The algorithm starts by constructing the *greedy graph*, which is the graph that takes the best incoming edge to each node, except the root. If the greedy graph contains a cycle, we *contract* the cycle into a single node c , and break the cycle by reweighting the enter edges, which are the edges that go into c .

¹⁹ Since computing the determinant can be done in $\mathcal{O}(N^3)$.

²⁰ The first two constraints are satisfied by the maximum-weight spanning tree, which we could compute with Kruskal’s algorithm. However, it does not adhere to the third constraint, because Kruskal’s algorithm only works for undirected graphs.



Then, we pick the greedy graph from the contracted graph. If there is more than one edge emanating from the root, we need to delete edges outgoing from the root to satisfy the single-root constraint. We remove the edge with the lowest *swap score*, which is the difference between the next-best incoming edge and the current incoming edge of each node in the graph.

If there is still a cycle, contract again, and continue doing this recursively until there are no more cycles. Then, expand all the cycles by picking the edges that are not canceled by the greedy graph with the contracted node.

Figure 9. Chu-Liu-Edmonds algorithm. The second graph is the greedy graph. The third graph is the contracted graph, where we contract the cycle into its own node. Then, we construct the greedy graph again. We have two edges emanating from the root, so we must eliminate the one with the lowest swap score. Then, we expand by choosing all the edges in the cycle that are not canceled.

12 Semantic parsing

Language cannot only be syntactically ambiguous, but also semantically ambiguous. Also, syntactically valid sentences do not necessarily mean anything. A good example of this is Chomsky's famous sentence "Colorless green dreams sleep furiously." In semantic analysis, we want to be able to reduce a natural language sentence, such as "Everyone loves someone else", to its logical form,

$$\forall p[\text{Person}(p) \rightarrow \exists q[\text{Person}(q) \wedge p \neq q \wedge \text{Loves}(p, q)]].$$

Note that this sentence is ambiguous, because we can swap the order of the quantifiers to a different logical form,

$$\exists p[\text{Person}(p) \rightarrow \forall q[\text{Person}(q) \wedge p \neq q \wedge \text{Loves}(q, p)]].$$

The challenge of semantic analysis, as always, is parsing a natural language sentence to its logical form. For this, we must use the *principle of compositionality*, which states that the meaning of a complex expression is a function of the meanings of that expression's constituent parts.²¹

12.1 Linear-indexed grammars

Linear-indexed grammars (LIG) [Aho, 1968] are a mildly context-sensitive family of grammars. The structure of an LIG is very similar to that of a CFG, but non-terminal symbols have an associated stack that can be passed to exactly one child,²²

$$N[\sigma] \rightarrow \alpha M[\sigma] \beta.$$

Definition 13 (Linear-indexed grammar). A linear-indexed grammar is a 5-tuple $\langle \mathcal{N}, S, I, \Sigma, \mathcal{R} \rangle$, such that

- \mathcal{N} is a set of non-terminal symbols, written as uppercase letters N, M, \dots ;
- $S \in \mathcal{N}$ is a distinguished start non-terminal;
- I is a finite set of indices, written as function letters f, g, h, \dots ;
- Σ is an alphabet of terminal symbols, written as lowercase letters a_1, a_2, \dots ;
- \mathcal{R} is a set of production rules of one of the following forms,

$$\begin{aligned} N[\sigma] &\rightarrow \alpha M[\sigma] \beta \\ N[\sigma] &\rightarrow \alpha M[f\sigma] \beta \\ N[f\sigma] &\rightarrow \alpha M[\sigma] \beta. \end{aligned}$$

²¹ This must hold, because otherwise we would not know the meaning of most sentences, since we have not heard most possible order of words that form a sentence.

Proverbs are edge cases to this principle, since those words have meaning together that are independent of its parts.

²² The fact that the stack can only be passed to exactly one child is what makes it linear.

Push

Pop

To make a grammar that encodes the subset $\{a^n b^n c^n \mid n \in \mathbb{N}\}$ is impossible with a CFG, but it is possible with an LIG, using the following grammar,

$$\begin{aligned} S[\sigma] &\rightarrow a S[f\sigma] c \\ S[\sigma] &\rightarrow T[\sigma] \\ T[f\sigma] &\rightarrow T[\sigma] b \\ T[] &\rightarrow \epsilon. \end{aligned}$$

Intuitively, it keeps track of how many times a and c are added on either side. Then, once it switches over to T from S after adding n times a and n times c on either side, we add n times b in the middle. The stack allows us to “count” to n .

12.2 Lambda calculus

Lambda calculus [Church, 1932] is a model for semantic analysis, based on the principle of compositionality. On a high level, lambda calculus first parses a sentence into simpler constituents, and then constructs the semantic representations using a bottom-up approach.



Figure 10. Lambda calculus on “Alex likes Bob.”

Lambda calculus has only the following rules,

- x, y, z, \dots are variables;
- $(\lambda x.f(x))$, which is the *lambda operator*, where x is a variable and $f(x)$ an expression;
- (MN) , which is an application of function M to an argument N , where they are both lambda terms.

A variable is bound if it belongs to a scope of abstraction holding its name. Otherwise, a variable is free. E.g., x is bound to λx and z is free in $((\lambda x.\lambda y.Likes(x, y))z)$.

To be able to represent natural language sentences, we also need the following,

- Constants that represent objects, denoted by *e.g.* ALEX, BOB, ...;
- Predicates that represent relations between objects, denoted by *e.g.* Teacher(\cdot), Likes(\cdot , \cdot), ...;
- Quantifiers \exists and \forall .

Furthermore, lambda calculus has two operations,

- α -conversion is the process of renaming the variable of a lambda operator and all of its bound occurrences. *E.g.*,

$$(\lambda x.x) \rightarrow (\lambda y.y);$$

- β -reduction is the process of applying one lambda term to another, *i.e.*, in $(\lambda x.MN)$, we replace all occurrences of x in M by N , and remove the lambda operator. *E.g.*,

$$(\lambda x.(\lambda y.xy)z) \rightarrow (\lambda y.zy).$$

However, sometimes a β -reduction would result in a free variable becoming bound. Then, we would first need to apply an α -conversion. Two lambda terms are equivalent if one can be obtained from the other after a series of α -conversions and β -reductions.

12.3 Combinatory logic

Combinatory logic [Curry et al., 1958] is an alternative to lambda calculus that formalizes the concept of computation and the construction of computable functions. Unlike lambda calculus, combinatory logic does not use abstractions. Instead, it uses complex functions using a few primitive higher order functions.

The basic terms of combinatory logic are

- x, y, z, \dots are variables;
- **I, S, K** are the primitive combinators, which are functions that map functions to functions.

Terms are then recursively constructed using the rule of application, where $(\mathbf{A} \mathbf{B})$ denotes applying **A** to **B**.

The following are the primitive combinators,²³

- $(\mathbf{K} x y) = x$ manufactures constant functions;
- $(\mathbf{S} x y z) = (x z (y z))$, which applies x to y after first substituting z into each of them;
- $(\mathbf{I} x) = x$ works as the identity function.²⁴

²³ In combinatory logic, parentheses are left-associative, *e.g.*, $(\mathbf{K} x y)$ means $((\mathbf{K} x) y)$ and $(\mathbf{S} x y z)$ means $((\mathbf{S} x) y z)$.

²⁴ **I** as a primitive combinator is not necessary, since it can be constructed from **S** and **K**,

$$\begin{aligned} ((\mathbf{S} \mathbf{K} \mathbf{K}) x) &= (\mathbf{S} \mathbf{K} \mathbf{K} x) \\ &= (\mathbf{K} x (\mathbf{K} x)) \\ &= x. \end{aligned}$$

Any lambda term is equivalent to a combinatory term that only uses the **S** and **K** combinators.

Further, there are the following combinators that are introduced for convenience,

$$(\mathbf{C} \ x \ y \ z) = ((x \ z) \ y)$$

Cross

$$(\mathbf{B} \ x \ y \ z) = (x \ (y \ z))$$

Composition

$$(\mathbf{T} \ x \ y) = (y \ z).$$

Type-raising

The **B** and **T** combinators will be used in combinatory categorial grammars.

12.4 Combinatory categorial grammars

Combinatory categorial grammars (CCG) are an efficiently parsable group of grammars that are mildly context-sensitive, which means that they have more expressive power than context-free grammars. CCGs allow us to model *coordination* and *cross-serial dependencies* in language, which CFGs cannot.

Definition 14 (Combinatory categorial grammar). A combinatory categorial grammar is a 5-tuple $\langle V_T, V_N, S, f, R \rangle$, such that

- V_T is a finite set of terminals;
- V_N is the finite set of atomic categories;
- $S \in V_N$ is the distinguished start category;
- f is a function mapping terminals $V_T \cup \{\epsilon\}$ to finite subsets of $C(V_N)$;
- R is a finite set of combinatory rules.

$C(V_N)$ is the infinite set of categories that contains all elements of V_N and recursively contains all elements such that if $c_1, c_2 \in C(V_N)$, then $c_1/c_2, c_1 \backslash c_2 \in C(V_N)$.

CCGs have two main parts: a lexicon that associates words with categories and rules that specify how categories can be combined into other categories. The lexicon contains all information specific to a given language, *i.e.*, valency, word order, and semantics. The structure information is encoded in the categories.²⁵

²⁵ Unlike CFGs, which encode structure in their rules.

The rules R are the following (inspired by combinatory logic),

$$\begin{array}{ll}
 X/Y \quad Y \implies X & (>) \\
 Y \quad X \backslash Y \implies X & (<) \\
 X/Y \quad Y/Z \implies X/Z & (\mathbf{B}_{>}) \\
 Y \backslash Z \quad X \backslash Y \implies X \backslash Z & (\mathbf{B}_{<}) \\
 \forall T \in V_N : X \implies T/(T \backslash X) & (\mathbf{T}_{>}) \\
 \forall T \in V_N : X \implies T \backslash (T/X), & (\mathbf{T}_{<})
 \end{array}$$

of which there are also generalized versions.

MARY	LIKES	JOHN
NP	$(S \backslash NP)/NP$	NP
MARY	$\lambda x. \lambda y. \text{Likes}(y, x)$	JOHN
$\xrightarrow{\quad}$		
$S \backslash NP$		
$\lambda y. \text{Likes}(y, \text{JOHN})$		
$\xrightarrow{\quad}$		
S		
$\text{Likes}(\text{MARY}, \text{JOHN})$		

(a) CCG derivation of “Mary likes John.”

WHAT	STATES	BORDER	TEXAS
$(S/(S \backslash NP))/N$	N	$(S \backslash NP)/NP$	NP
$\lambda f. \lambda g. \lambda x. f(x) \wedge g(x)$	$\lambda x. \text{State}(x)$	$\lambda x. \lambda y. \text{Borders}(y, x)$	TEXAS
$\xrightarrow{\quad}$		$\xrightarrow{\quad}$	
$S/(S \backslash NP)$		$S \backslash NP$	
$\lambda g. \lambda x. \text{State}(x) \wedge g(x)$		$\lambda y. \text{Borders}(y, \text{TEXAS})$	
$\xrightarrow{\quad}$		$\xrightarrow{\quad}$	
S		S	
$\lambda x. \text{State}(x) \wedge \text{Borders}(x, \text{TEXAS})$			

(b) CCG derivation of “What states border Texas?”

Figure 11. Simple example CCG derivations.

13 Transformers

Attention is a mechanism in neural networks that a model can learn to make predictions by selectively attending to a given set of data by using query q , key k , and value v vector representations. The query and key vectors are used to determine how much weight should be given to the value vector.²⁶ The weights are computed by $a_{ij} = \text{softmax}(q_j^\top k_i)$, so the values after the attention block can be computed as follows,

$$\text{att}(x_i) = \sum_j a_{ij} v_j,$$

which is a linear combination of the values according to the attention weights computed by the query and keys.

Self-attention blocks learn the query, key, and value representations from data. More specifically, it learns matrices W_Q , W_K , and W_V and computes the vectors from these matrices:

$$\begin{aligned} q_i &= W_Q^\top x_i \\ k_i &= W_K^\top x_i \\ v_i &= W_V^\top x_i. \end{aligned}$$

Then, we can use these to compute the output of the self-attention block:

$$\text{self-att}(X) = \text{softmax}\left(\frac{(W_Q^\top X)^\top (W_K^\top X)}{\sqrt{d_q}}\right) W_V^\top X,$$

where d_q is the square root of the dimensionality of the query and key vectors. Furthermore, we need to add a positional encoding to provide ordering information to the model.²⁷ This is done by a sinusoidal positional encoding and are simply combined with X by addition.

Transformers [Vaswani et al., 2017] use multi-headed self-attention, which is a module where self-attention is applied M times independently to the data. Thus, this module learns M different ways of looking at the same dataset. The outputs of each self-attention block is concatenated and linearly transformed to the expected dimensionality. Transformers follow this by normalization and MLP layers, as can be seen in Figure 13.

13.1 Translation

Translation is a sequence-to-sequence problem, where we want to compute the probability that y is the translation of x . We can do this with transformers by encoding the input sequence x using encoders, and feeding this representation of the input to a decoder. The decoder takes as input the input sequence and the already generated (incomplete) sequence $y_{<i}$. It then runs the encoder as in Figure 13 and projects the output to a probability distribution over tokens using a linear layer, followed by softmax.

²⁶ Note the parallel with dictionaries/hashmaps in programming languages, but, in the attention mechanism, we do a “soft-lookup”.



Figure 12. Self-attention mechanism.

²⁷ The self-attention operation is permutation equivariant.



Figure 13. Transformer encoder architecture.

Figure 14. Translation architecture

In previous problem statements, we were usually able to compute the globally maximum output using some dynamic program. However, in this case, that is not possible. This model is only locally normalized. But, we can do something more optimal than greedily picking the next token. We can view this problem as a graph search problem over possible output sentences. We still cannot explore every path of the $\mathcal{O}(|\Sigma|^{N_{\max}})$ paths, but we use beam search. Beam search keeps track of a maximum of k paths at any given time. Then, it prunes paths that are the worst. This is still not globally optimal, but it is a bit better than the greedy approach.

References

- Alfred V Aho. Indexed grammars—an extension of context-free grammars. *Journal of the ACM (JACM)*, 15(4):647–671, 1968.
- Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. A neural probabilistic language model. *Advances in neural information processing systems*, 13, 2000.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation, 2014.
- Yoeng-Jin Chu. On the shortest arborescence of a directed graph. *Scientia Sinica*, 14:1396–1400, 1965.
- Alonzo Church. A set of postulates for the foundation of logic. *Annals of mathematics*, pages 346–366, 1932.
- John Cocke. *Programming languages and their compilers: Preliminary notes*. New York University, 1969.
- Haskell Brooks Curry, Robert Feys, William Craig, J Roger Hindley, and Jonathan P Seldin. *Combinatory logic*, volume 1. North-Holland Amsterdam, 1958.
- Jack Edmonds et al. Optimum branchings. *Journal of Research of the national Bureau of Standards B*, 71(4):233–240, 1967.
- Jeffrey L Elman. Finding structure in time. *Cognitive science*, 14(2):179–211, 1990.
- Robert W Floyd. Algorithm 97: shortest path. *Communications of the ACM*, 5(6):345, 1962.
- Simon Haykin. *Neural networks: a comprehensive foundation*. Prentice Hall PTR, 1994.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Mohit Iyyer, Varun Manjunatha, Jordan Boyd-Graber, and Hal Daumé III. Deep unordered composition rivals syntactic methods for text classification. 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)*, pages 1681–1691, 2015.
- Tadao Kasami. An efficient recognition and syntax-analysis algorithm for context-free languages. *Coordinated Science Laboratory Report no. R-257*, 1966.

Gustav Kirchhoff. Ueber die auflösung der gleichungen, auf welche man bei der untersuchung der linearen vertheilung galvanischer ströme geführt wird. *Annalen der Physik*, 148(12):497–508, 1847.

Terry Koo, Amir Globerson, Xavier Carreras Pérez, and Michael Collins. Structured prediction models via the matrix-tree theorem. In *Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 141–150, 2007.

Daniel J Lehmann. Algebraic structures for transitive closure. *Theoretical Computer Science*, 4(1):59–76, 1977.

Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

Bo Pang, Lillian Lee, et al. Opinion mining and sentiment analysis. *Foundations and Trends® in information retrieval*, 2(1–2):1–135, 2008.

Robert Endre Tarjan. Finding optimum branchings. *Networks*, 7(1):25–35, 1977.

William T Tutte. The dissection of equilateral triangles into equilateral triangles. In *Mathematical Proceedings of the Cambridge Philosophical Society*, volume 44, pages 463–482. Cambridge University Press, 1948.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

Stephen Warshall. A theorem on boolean matrices. *Journal of the ACM (JACM)*, 9(1):11–12, 1962.

Daniel H Younger. Recognition and parsing of context-free languages in time n^3 . *Information and control*, 10(2):189–208, 1967.