



Formal Aspects of Language Modeling

Ryan Cotterell, Anej Svete, Clara Meister,
Tianyu Liu and Li Du

Thursday 9th November, 2023

Contents

1	Introduction	5
1.1	Introduction	5
2	Probabilistic Foundations	7
2.1	An Invitation to Language Modeling	7
2.2	A Measure-theoretic Foundation	10
2.3	Language Models: Distributions over Strings	13
2.3.1	Sets of Strings	13
2.3.2	Defining a Language Model	14
2.4	Global and Local Normalization	16
2.4.1	Globally Normalized Language Models	16
2.4.2	Locally Normalized Language Models	17
2.5	Tight Language Models	22
2.5.1	Tightness	22
2.5.2	Defining the probability measure of an LNM	24
2.5.3	Interpreting the Constructed Probability Space	29
2.5.4	Characterizing Tightness	30
3	Modeling Foundations	37
3.1	Representation-based Language Models	38
3.1.1	Vector Space Representations	38
3.1.2	Compatibility of Symbol and Context	42
3.1.3	Projecting onto the Simplex	43
3.1.4	Representation-based Locally Normalized Models	47
3.1.5	Tightness of Softmax Representation-based Models	47
3.2	Estimating a Language Model from Data	49
3.2.1	Data	49
3.2.2	Language Modeling Objectives	49
3.2.3	Parameter Estimation	54
3.2.4	Regularization Techniques	57
4	Classical Language Models	59
4.1	Finite-state Language Models	59
4.1.1	Weighted Finite-state Automata	59
4.1.2	Finite-state Language Models	66
4.1.3	Normalizing Finite-state Language Models	67
4.1.4	Tightness of Finite-state Models	72
4.1.5	The n -gram Assumption and Subregularity	75
4.1.6	Representation-based n -gram Models	78
4.2	Pushdown Language Models	83
4.2.1	Human Language is not Finite-state	83
4.2.2	Context-free Grammars	84
4.2.3	Weighted Context-free Grammars	88

4.2.4	Context-free Language Models	91
4.2.5	Tightness of Context-free Language Models	92
4.2.6	Normalizing Weighted Context-free Grammars	94
4.2.7	Pushdown Automata	96
4.2.8	Pushdown Language Models	100
4.2.9	Multi-stack Pushdown Automata	100
4.3	Exercises	103
5	Neural Network Language Models	105
5.1	Recurrent Neural Language Models	106
5.1.1	Human Language is Not Context-free	106
5.1.2	Recurrent Neural Networks	107
5.1.3	General Results on Tightness	112
5.1.4	Elman and Jordan Networks	114
5.1.5	The Vanishing and Exploding Gradient Problem	116
5.1.6	Variations on Recurrent Networks	120
5.1.7	Computational Expressiveness of Recurrent Neural Networks	124
5.2	Transformer-based Language Models	163
5.2.1	Informal Motivation of the Transformer Architecture	163
5.2.2	A Formal Definition of Transformers	164
5.2.3	Tightness of Transformer-based Language Models	174
5.3	Computational Expressiveness of Transformers	175

Chapter 1

Introduction

1.1 Introduction

Welcome to the class notes for the first half of Large Language Models (263-5354-00L). The course comprises an omnibus introduction to language modeling. The first half of the lectures focuses on a formal treatment of the subject. The second half focuses on the practical aspects of implementing a language model and its applications. Many universities are offering similar courses at the moment, e.g., CS324 at Stanford University (<https://stanford-cs324.github.io/winter2022/>) and CS 600.471 (<https://self-supervised.cs.jhu.edu/sp2023/>) at Johns Hopkins University. Their syllabi may serve as useful references.

Chapter 2

Probabilistic Foundations

2.1 An Invitation to Language Modeling

The first module of the course focuses on *defining* a language model mathematically. To see why such a definition is nuanced, we are going to give an informal definition of a language model and demonstrate two ways in which that definition breaks and fails to meet our desired criteria.

Definition 2.1.1 (Informal). *Given an alphabet¹ Σ and a distinguished end-of-sequence symbol $\text{EOS} \notin \Sigma$, a language model is a collection of conditional probability distributions $p(y \mid \mathbf{y})$ for $y \in \Sigma \cup \{\text{EOS}\}$ and $\mathbf{y} \in \Sigma^*$, the Kleene closure of Σ .² $p(y \mid \mathbf{y})$ therefore represents the probability of y being the next token given the history \mathbf{y} .*

Definition 2.1.1 is the definition of a language model that is implicitly assumed in most papers on language modeling. We say implicitly since most technical papers on language modeling simply write down the following autoregressive factorization

$$p(\mathbf{y}) = p(y_1 \cdots y_T) = p(\text{EOS} \mid \mathbf{y}) \prod_{t=1}^T p(y_t \mid \mathbf{y}_{<t}) \quad (2.1)$$

as the probability of a string according to the distribution p .³ The part that is left implicit in Eq. (2.1) is whether or not p is indeed a probability distribution and, if it is, over what space. The natural assumption in Definition 2.1.1 is that p is a distribution over Σ^* , i.e., the set of all *finite* strings⁴ over an alphabet Σ . However, in general, it is not true that all such collections of conditionals will yield a valid probability distribution over Σ^* ; some may “leak” probability mass to infinite sequences.⁵ More subtly, we additionally have to be very careful when dealing with uncountably infinite spaces lest we run into a classic paradox, which we will demonstrate below.

We highlight the two issues above with two very simple examples. The first example is a well-known paradox in probability theory.

Example 2.1.1 (Infinite Coin Toss). *Consider the infinite independent fair coin toss model, where we aim to place a distribution over $\{\text{H}, \text{T}\}^\infty$, the (uncountable) set of infinite sequences of $\{\text{H}, \text{T}\}$ (H represents the event of throwing heads and T the event of throwing tails). Intuitively, such a distribution corresponds to a “language model” as defined above in which for all $\mathbf{y}_{<t}$, $p(\text{H} \mid \mathbf{y}_{<t}) = p(\text{T} \mid \mathbf{y}_{<t}) = \frac{1}{2}$ and $p(\text{EOS} \mid \mathbf{y}_{<t}) = 0$. However, each individual infinite sequence over $\{\text{H}, \text{T}\}$ should also be assigned probability $(\frac{1}{2})^\infty = 0$. Without*

¹An alphabet is a *finite*, non-empty set.

²We will define the Kleene closure in more detail later in the course.

³Many authors (erroneously) avoid writing EOS for concision.

⁴Some authors assert that strings are by definition finite.

⁵However, the converse *is* true: All valid distributions over Σ^* may be factorized as the above.

a formal foundation, one arrives at the following paradox:

$$\begin{aligned} 1 &= p(\{\text{H}, \text{T}\}^\infty) = p\left(\bigcup_{\omega \in \{\text{H}, \text{T}\}^\infty} \{\omega\}\right) \\ &= \sum_{\omega \in \{\text{H}, \text{T}\}^\infty} p(\{\omega\}) = \sum_{\omega \in \{\text{H}, \text{T}\}^\infty} 0 \stackrel{?}{=} 0. \end{aligned}$$

The second example is more specific to language modeling. As we stated above, an implicit assumption made by most language modeling papers is that a language model constitutes a distribution over Σ^* . However, in our next example, we show that a collection of conditions that satisfy Definition 2.1.1 may not sum to 1 if the sum is restricted to elements of Σ^* . This means that it is not a priori clear what space our probability distribution is defined over. ⁶

Example 2.1.2 (Possibly Finite Coin Toss). *Consider now the possibly finite “coin toss” model with a rather peculiar coin: when tossing the coin for the first time, both H and T are equally likely. After the first toss, however, the coin gets stuck: If $y_1 = \text{H}$, we can only ever toss another H again, whereas if $y_1 = \text{T}$, the next toss can result in another T or “end” the sequence of throws (EOS) with equal probability. We, therefore, model a probability distribution over $\{\text{H}, \text{T}\}^* \cup \{\text{H}, \text{T}\}^\infty$, the set of finite and infinite sequences of tosses. Formally:⁷*

$$\begin{aligned} p(\text{H} \mid \mathbf{y}_{<1}) &= p(\text{T} \mid \mathbf{y}_{<1}) = \frac{1}{2} \\ p(\text{H} \mid \mathbf{y}_{<t}) &= \begin{cases} 1 & \text{if } t > 1 \text{ and } y_{t-1} = \text{H} \\ 0 & \text{if } t > 1 \text{ and } y_{t-1} = \text{T} \end{cases} \\ p(\text{T} \mid \mathbf{y}_{<t}) &= \begin{cases} \frac{1}{2} & \text{if } t > 1 \text{ and } y_{t-1} = \text{T} \\ 0 & \text{if } t > 1 \text{ and } y_{t-1} = \text{H} \end{cases} \\ p(\text{EOS} \mid \mathbf{y}_{<t}) &= \begin{cases} \frac{1}{2} & \text{if } t > 1 \text{ and } y_{t-1} = \text{T} \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

If you are familiar with (weighted) finite-state automata,⁸ you can imagine the model as depicted in Fig. 2.1.

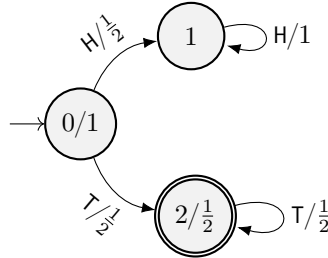


Figure 2.1: Graphical depiction of the possibly finite coin toss model. The final weight $\frac{1}{2}$ of the state 2 corresponds to the probability $p(\text{EOS} \mid y_{t-1} = \text{T}) = \frac{1}{2}$.

It is easy to see that this model only places the probability of $\frac{1}{2}$ on finite sequences of tosses. If we were only interested in those (analogously to how we are only interested in finite strings when modeling language), yet still allowed the model to specify the probabilities as in this example, the resulting probability distribution would not model what we require.

It takes some mathematical heft to define a language model in a manner that avoids such paradoxes. The tool of choice for mathematicians is measure theory, as it allows us to define probability over uncountable

⁶This also holds for the first example.

⁷Note that $p(\text{H} \mid \mathbf{y}_{<1}) = p(\text{H} \mid \varepsilon)$ and $p(\text{T} \mid \mathbf{y}_{<1}) = p(\text{T} \mid \varepsilon)$.

⁸They will be formally introduced in §4.1.5

sets⁹ in a principled way. Thus, we begin our formal treatment of language modeling with a primer of measure theory in §2.2. Then, we will use concepts discussed in the primer to work up to a formal definition of a language model.

⁹As stated earlier, $\{\mathsf{H}, \mathsf{T}\}^\infty$ is uncountable. It's easy to see there exists a surjection from $\{\mathsf{H}, \mathsf{T}\}^\infty$ to the binary expansion of the real interval $(0, 1]$. Readers who are interested in more details and mathematical implications can refer to §1 in [Billingsley \(1995\)](#).

2.2 A Measure-theoretic Foundation

At their core, (large) language models are an attempt to place a probabilistic distribution over natural language utterances. However, our toy examples in Example 2.1.1 and ?? in the previous section reveal that it is quite nuanced to get a satisfying definition of a language model. Thus, our first step forward is to review the basics of rigorous probability theory,¹⁰ the tools we need to come to a satisfying definition. Our course will assume that you have had some exposure to rigorous probability theory before, and just review the basics. However, it is also possible to learn the basics of rigorous probability on the fly during the course if it is new to you. Specifically, we will cover *measure-theoretic* foundations of probability theory. This might come as a bit of a surprise since we are mostly going to be talking about discrete objects—namely, strings. However, as we will see in §2.5 soon, formal treatment of language modeling indeed requires some mathematical rigor from measure theory.

The goal of measure-theoretic probability is to assign probabilities to *subsets* of an **outcome space** Ω . However, in the course of the study of measure theory, it has become clear that for many common Ω , it is impossible to assign probabilities in a way that satisfies a set of reasonable desiderata.¹¹ Consequently, the standard approach to probability theory resorts to only assigning probability to certain “nice” (but not all) subsets of Ω , which are referred to as **events** or **measurable subsets**, as in the theory of integration or functional analysis. The set of measurable subsets is commonly denoted as \mathcal{F} (Definition 2.2.1) and a probability measure $\mathbb{P} : \mathcal{F} \rightarrow [0, 1]$ is the function that assigns a probability to each measurable subset. The triple $(\Omega, \mathcal{F}, \mathbb{P})$ is collectively known as a probability space (Definition 2.2.2). As it turns out, the following simple and reasonable requirements imposed on \mathcal{F} and \mathbb{P} are enough to rigorously discuss probability.

Definition 2.2.1. Let $\mathcal{P}(\Omega)$ be the power set of Ω . Then $\mathcal{F} \subseteq \mathcal{P}(\Omega)$ is called a **σ -algebra** (or **σ -field**) over Ω if the following conditions hold:

- 1) $\Omega \in \mathcal{F}$,
- 2) if $\mathcal{E} \in \mathcal{F}$, then $\mathcal{E}^c \in \mathcal{F}$,
- 3) if $\mathcal{E}_1, \mathcal{E}_2, \dots$ is a finite or infinite sequence of sets in \mathcal{F} , then $\bigcup_n \mathcal{E}_n \in \mathcal{F}$.

If \mathcal{F} is a σ -algebra over Ω , we call the tuple (Ω, \mathcal{F}) a **measurable space**.

Example 2.2.1 (σ -algebras). Let Ω be any set. Importantly, there is no one way to construct a σ -algebra over Ω :

1. The family consisting of only the empty set \emptyset and the set Ω , i.e., $\mathcal{F} \stackrel{\text{def}}{=} \{\emptyset, \Omega\}$, is called the **minimal** or **trivial**.
2. The full power set $\mathcal{F} \stackrel{\text{def}}{=} \mathcal{P}(\Omega)$ is called the **discrete σ -algebra**.
3. Given $\mathcal{A} \subseteq \Omega$, the family $\mathcal{F} \stackrel{\text{def}}{=} \{\emptyset, \mathcal{A}, \Omega \setminus \mathcal{A}, \Omega\}$ is a σ -algebra induced by \mathcal{A} .
4. Suppose we are rolling a six-sided die. There are six events that can happen: we can roll any of the numbers 1–6. In this case, we will then define the set of outcomes Ω as $\Omega \stackrel{\text{def}}{=} \{\text{The number observed is } n \mid n = 1, \dots, 6\}$. There are of course multiple ways to define an event space \mathcal{F} and with it a σ -algebra over this sample space. By definition, $\emptyset \in \mathcal{F}$ and $\Omega \in \mathcal{F}$. One way to intuitively construct a σ -algebra is to consider that all individual events (observing any number) are possible, meaning that we would like to later assign probabilities to them (see Definition 2.2.2). This means that we should include individual singleton events in the event space: $\{\text{The number observed is } n\} \in \mathcal{F}$ for $n = 1, \dots, 6$. It is easy to see that in this case, to satisfy the axioms in Definition 2.2.1, the resulting event space should be $\mathcal{F} = \mathcal{P}(\Omega)$.

You might want to confirm these are indeed σ -algebras by checking them against the axioms in Definition 2.2.1.

A measurable space guarantees that operations on countably many sets are always valid, and hence permits the following definition.

¹⁰By rigorous probability theory we mean a measure-theoretic treatment of probability theory.

¹¹Measure theory texts commonly discuss such desiderata and the dilemma that comes with it. See, e.g., Chapter 7 in Tao (2016), Chapter 3 in Royden (1988) or Chapter 3 in Billingsley (1995). We also give an example later.

Definition 2.2.2. A **probability measure** \mathbb{P} over a measurable space (Ω, \mathcal{F}) is a function $\mathbb{P} : \mathcal{F} \rightarrow [0, 1]$ such that

- 1) $\mathbb{P}(\Omega) = 1$,
- 2) if $\mathcal{E}_1, \mathcal{E}_2, \dots$ is a countable sequence of disjoint sets in \mathcal{F} , then $\mathbb{P}(\bigcup_n \mathcal{E}_n) = \sum_n \mathbb{P}(\mathcal{E}_n)$.

In this case we call $(\Omega, \mathcal{F}, \mathbb{P})$ a **probability space**.

As mentioned, measure-theoretic probability only assigns probabilities to “nice” subsets of Ω . In fact, it is often impossible to assign a probability measure to every single subset of Ω and we must restrict our probability space to a strict subset of $\mathcal{P}(\Omega)$. More precisely, the sets $\mathcal{B} \subseteq \Omega$ for which a probability (or more generally, a *volume*) can not be defined are called *non-measurable sets*. An example of under certain assumptions such sets is the Vitali set¹² See also Appendix A.2 in Durrett (2019).

Later, we will be interested in modeling probability spaces over sets of (infinite) sequences. By virtue of a theorem due to Carathéodory, there is a natural way to construct such a probability space for sequences (and many other spaces) that behaves in accordance with our intuition, as we will clarify later. Here, we shall lay out a few other necessary definitions.

Definition 2.2.3. $\mathcal{A} \subseteq \mathcal{P}(\Omega)$ is called an **algebra** (or *field*) over Ω if

- 1) $\Omega \in \mathcal{A}$,
- 2) if $\mathcal{E} \in \mathcal{A}$, then $\mathcal{E}^c \in \mathcal{A}$,
- 3) if $\mathcal{E}_1, \mathcal{E}_2 \in \mathcal{A}$, then $\mathcal{E}_1 \cup \mathcal{E}_2 \in \mathcal{A}$.

Definition 2.2.4. Let \mathcal{A} be an algebra over some set Ω . A **probability pre-measure** over (Ω, \mathcal{A}) is a function $\mathbb{P}_0 : \mathcal{A} \rightarrow [0, 1]$ such that

- 1) $\mathbb{P}_0(\Omega) = 1$,
- 2) if $\mathcal{E}_1, \mathcal{E}_2, \dots$ is a (countable) sequence of disjoint sets in \mathcal{A} whose (countable) union is also in \mathcal{A} , then $\mathbb{P}_0(\bigcup_{n=1}^{\infty} \mathcal{E}_n) = \sum_{n=1}^{\infty} \mathbb{P}_0(\mathcal{E}_n)$.

Note that the only difference between a σ -algebra (Definition 2.2.1) and an algebra is that condition 3 is weakened from countable to finite, and the only difference between a probability measure (Definition 2.2.2) and a pre-measure is that the latter is defined with respect to an algebra instead of a σ -algebra.

The idea behind Carathéodory’s extension theorem is that there is often a simple construction of an algebra \mathcal{A} over Ω such that there is a natural way to define a probability pre-measure. One can then *extend* this probability pre-measure to a probability measure that is both minimal and unique in a precise sense. For example, the standard Lebesgue measure over the real line can be constructed this way.

Finally, we define random variables.

Definition 2.2.5 (Random variable). A mapping $x : \Omega \rightarrow \mathcal{S}$ between two measurable spaces (Ω, \mathcal{F}) and $(\mathcal{S}, \mathcal{T})$ is an $(\mathcal{S}, \mathcal{T})$ -valued **random variable**, or a *measurable mapping*, if, for all $\mathcal{B} \in \mathcal{T}$,

$$x^{-1}(\mathcal{B}) \stackrel{\text{def}}{=} \{\omega \in \Omega : x(\omega) \in \mathcal{B}\} \in \mathcal{F}. \quad (2.2)$$

Any measurable function (random variable) induces a new probability measure on the *output* σ -algebra based on the one defined on the original σ -algebra. This is called the **pushforward measure** (cf. §2.4 in Tao, 2011), which we will denote by \mathbb{P}_* , given by

$$\mathbb{P}_*(x \in \mathcal{E}) \stackrel{\text{def}}{=} \mathbb{P}(x^{-1}(\mathcal{E})), \quad (2.3)$$

that is, the probability of the result of x being in some event \mathcal{E} is determined by the probability of the event of all the elements which x maps into \mathcal{E} , i.e., the pre-image of \mathcal{E} given by x .

¹²See https://en.wikipedia.org/wiki/Non-measurable_set and https://en.wikipedia.org/wiki/Vitali_set.

Example 2.2.2 (Random Variables). *We give some simple examples of random variables.*

1. Let Ω be the set of possible outcomes of throwing a fair coin, i.e., $\Omega \stackrel{\text{def}}{=} \{\text{T}, \text{H}\}$. Define $\mathcal{F} \stackrel{\text{def}}{=} \mathcal{P}(\Omega)$, $\mathcal{S} \stackrel{\text{def}}{=} \{0, 1\}$, and $\mathcal{T} \stackrel{\text{def}}{=} \mathcal{P}(\mathcal{S})$. Then, the random variable

$$\mathbf{x} : \begin{cases} \text{T} \mapsto 0 \\ \text{H} \mapsto 1 \end{cases}$$

assigns tails (T) the value 0 and heads (H) the value 1.

2. Consider the probability space of throwing two dice (similar to Example 2.2.1) where $\Omega = \{(i, j) : i, j = 1, \dots, 6\}$ where the element (i, j) refers to rolling i on the first and j on the second die and $\mathcal{F} = \mathcal{P}(\Omega)$. Define $\mathcal{S} \stackrel{\text{def}}{=} \mathbb{Z}$ and $\mathcal{T} \stackrel{\text{def}}{=} \mathcal{P}(\mathcal{S})$. Then, the random variable

$$\mathbf{x} : (i, j) \mapsto i + j$$

is an $(\mathcal{S}, \mathcal{T})$ -valued random variable which represents the sum of two dice.

2.3 Language Models: Distributions over Strings

Language models are defined as probability distributions over sequences of words, referred to as utterances. This chapter delves into the formalization of the term “utterance” and introduces fundamental concepts such as the alphabet, string, and language. Utilizing these concepts, a formal definition of a language model is presented, along with a discussion on the intricacies of defining distributions over infinite sets.

2.3.1 Sets of Strings

We begin by defining the very basic notions of alphabets and strings, where we take inspiration from **formal language theory**. First and foremost, formal language theory concerns itself with *sets of structures*. The simplest structure it considers is a **string**. So what is a string? We start with the notion of an alphabet.

Definition 2.3.1 (Alphabet). An **alphabet** is a finite, non-empty set. In this course, we will denote an alphabet using Greek capital letters, e.g., Σ and Δ . We refer to the elements of an alphabet as **symbols** or **letters** and will denote them with lowercase letters: a, b, c .

Definition 2.3.2 (String). A **string**¹³ over an alphabet is any finite sequence of letters. Strings made up of symbols from Σ will be denoted by bolded Latin letters, e.g., $\mathbf{y} = y_1 \cdots y_T$ where each $y_n \in \Sigma$.

The length of a string, written as $|\mathbf{y}|$, is the number of letters it contains. Usually, we will use T to denote $|\mathbf{y}|$ more concisely whenever the usage is clear from the context. There is only one string of length zero, which we denote with the distinguished symbol ε and refer to as the *empty string*. By convention, ε is *not* an element of the original alphabet.

New strings are formed from other strings and symbols with **concatenation**. Concatenation, denoted with $\mathbf{x} \circ \mathbf{y}$ or just \mathbf{xy} , is an associative operation on strings. Formally, the concatenation of two words \mathbf{y} and \mathbf{x} is the word $\mathbf{y} \circ \mathbf{x} = \mathbf{yx}$, which is obtained by writing the second argument after the first one. The result of concatenating with ε from either side results in the original string, which means that ε is the **unit** of concatenation and the set of all words over an alphabet with the operation of concatenation forms a **monoid**.

We have so far only defined strings as individual sequences of symbols. To give our strings made up of symbols in Σ a set to live in, we now define Kleene closure of an alphabet Σ .

Definition 2.3.3 (Kleene Star). Let Σ be an alphabet. The **Kleene star** Σ^* is defined as

$$\Sigma^* = \bigcup_{n=0}^{\infty} \Sigma^n \quad (2.4)$$

where

$$\Sigma^n \stackrel{\text{def}}{=} \underbrace{\Sigma \times \cdots \times \Sigma}_{n \text{ times}} \quad (2.5)$$

Note that we define $\Sigma^0 \stackrel{\text{def}}{=} \{\varepsilon\}$. We call the Σ^* the **Kleene closure** of the alphabet Σ . We also define

$$\Sigma^+ \stackrel{\text{def}}{=} \bigcup_{n=1}^{\infty} \Sigma^n = \Sigma \Sigma^*. \quad (2.6)$$

Finally, we also define the set of all infinite sequences of symbols from some alphabet Σ as Σ^∞ .

Definition 2.3.4. Let Σ be an alphabet. The set of all **infinite sequences** over Σ is defined as:

$$\Sigma^\infty \stackrel{\text{def}}{=} \underbrace{\Sigma \times \cdots \times \Sigma}_{\infty\text{-times}}, \quad (2.7)$$

Since strings are canonically *finite* in computer science, we will explicitly use the terms infinite sequence or infinite string to refer to elements of Σ^∞ .

¹³A string is also referred to as a **word**, which continues with the linguistic terminology.

More informally, we can think of Σ^* as the set which contains ε and all (finite-length) strings which can be constructed by concatenating arbitrary symbols from Σ . Σ^+ , on the other hand, does *not* contain ε , but contains all other strings of symbols from Σ . The Kleene closure of an alphabet is a *countably infinite* set (this will come into play later!). In contrast, the set Σ^∞ is *uncountably infinite* for any Σ such that $|\Sigma| \geq 2$.

The notion of the Kleene closure leads us very naturally to our next definition.

Definition 2.3.5 (Formal language). *Let Σ be an alphabet. A **language** L is a subset of Σ^* .*

That is, a language is just a specified subset of all possible strings made up of the symbols in the alphabet. This subset can be specified by simply enumerating a finite set of strings, or by a *formal model*. We will see examples of those later. Importantly, these strings are *finite*. If not specified explicitly, we will often assume that $L = \Sigma^*$.

A note on terminology. As we mentioned, these definitions are inspired by formal language theory. We defined strings as our main structures of interest and symbols as their building blocks. When we talk about natural language, the terminology is often slightly different: we may refer to the basic building blocks (symbols) as **tokens** or **words** (which might be composed of one or more *characters* and form some form of “words”) and their compositions (strings) as **sequences** or **sentences**. Furthermore, what we refer to here as an alphabet may be called a **vocabulary** (of words or tokens) in the context of natural language. Sentences are therefore concatenations of words from a vocabulary in the same way that strings are concatenations of symbols from an alphabet.

Example 2.3.1 (Kleene Closure). *Let $\Sigma = \{a, b, c\}$. Then*

$$\Sigma^* = \{\varepsilon, a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, aab, aac, \dots\}.$$

Examples of a languages over this alphabet include $L_1 \stackrel{\text{def}}{=} \{a, b, ab, ba\}$, $L_2 \stackrel{\text{def}}{=} \{\mathbf{y} \in \Sigma^ \mid y_1 = a\}$, and $L_3 \stackrel{\text{def}}{=} \{\mathbf{y} \in \Sigma^* \mid |\mathbf{y}| \text{ is even}\}$.*

Next, we introduce two notions of subelements of strings.

Definition 2.3.6 (String Subelements). *A **subsequence** of a string \mathbf{y} is defined as a sequence that can be formed from \mathbf{y} by deleting some or no symbols, leaving the order untouched. A **substring** is a contiguous subsequence. For instance, ab and bc are substrings and subsequences of $\mathbf{y} = abc$, while ac is a subsequence but not a substring. **Prefixes** and **suffixes** are special cases of substrings. A prefix is a substring of \mathbf{y} that shares the same first letter as \mathbf{y} and a suffix is a substring of \mathbf{y} that shares the same last letter as \mathbf{y} . We will also denote a prefix $y_1 \dots y_{n-1}$ of the string $\mathbf{y} = y_1 \dots y_T$ as $\mathbf{y}_{<n}$. We will also use the notation $\mathbf{y} \triangleleft \mathbf{y}'$ to denote that \mathbf{y} is a suffix of \mathbf{y}' .*

2.3.2 Defining a Language Model

We are now ready to introduce the main interest of the entire lecture series: language models.

Definition 2.3.7 (Language model). *Let Σ be an alphabet. A **language model** is a (discrete) distribution p_{LM} over Σ^* .*

Example 2.3.2 (A very simple language model). *Let $\Sigma \stackrel{\text{def}}{=} \{a\}$. For $n \in \mathbb{N}_{\geq 0}$, define*

$$p_{LM}(a^n) \stackrel{\text{def}}{=} 2^{-(n+1)},$$

where $a^0 \stackrel{\text{def}}{=} \varepsilon$ and $a^n \stackrel{\text{def}}{=} \underbrace{a \dots a}_{n \text{ times}}$.

We claim that p_{LM} is a language model. To see that, we verify that it is a valid probability distribution over Σ^ . It is easy to see that $p_{LM}(a^n) \geq 0$ for any n . Additionally, we see that the probabilities of finite sequences indeed sum to 1:*

$$\sum_{\mathbf{y} \in \Sigma^*} p_{LM}(\mathbf{y}) = \sum_{n=0}^{\infty} p_{LM}(a^n) = \sum_{n=0}^{\infty} 2^{-(n+1)} = \frac{1}{2} \sum_{n=0}^{\infty} 2^{-n} = \frac{1}{2} \frac{1}{1 - \frac{1}{2}} = 1.$$

In our formal analysis of language models, we will also often refer to the *language* defined by a language model.

Definition 2.3.8. Let p_{LM} be a language model. The **weighted language** of p_{LM} is defined as

$$L(p_{LM}) \stackrel{\text{def}}{=} \{(\mathbf{y}, p_{LM}(\mathbf{y})) \mid \mathbf{y} \in \Sigma^*\} \quad (2.8)$$

Example 2.3.3. The language of the language model from Example 2.3.2 is

$$L(p_{LM}) \stackrel{\text{def}}{=} \left\{ \left(a^n, 2^{-(n+1)} \right) \mid n \in \mathbb{N}_{\geq 0} \right\} \quad (2.9)$$

A language model is itself a very simple concept—it is simply a distribution that weights strings (natural utterances) by their probabilities to occur in a particular language. Note that we have not said anything about how we can represent or model this distribution yet. Besides, for any (natural) language, the ground-truth language model p_{LM} is of course *unknown* and complex. The next chapter, therefore, discusses in depth the computational models which we can use to try to tractably represent distributions over strings and ways of *approximating* (learning) the ground-truth distribution based on finite datasets using such models.

2.4 Global and Local Normalization

The previous chapter introduced a formal definition of a language as a set of strings and the definition of a language model as a distribution over strings. We now delve into a potpourri of technical questions to complete the theoretical minimum for discussing language models. While doing so, we will introduce (and begin to answer) three fundamental questions in the first part of the course. We will introduce them later in the section.

A note on terminology. Unfortunately, we will encounter some ambiguous terminology. In §2.5, we explicitly define a language model as a valid probability distribution over Σ^* , the Kleene closure of some alphabet Σ , which means that $\sum_{\mathbf{y} \in \Sigma^*} p_{\text{LM}}(\mathbf{y}) = 1$. As we will see later, this means that the model is *tight*, whereas it is *non-tight* if $\sum_{\mathbf{y} \in \Sigma^*} p_{\text{LM}}(\mathbf{y}) < 1$. Definitionally, then, all language models are tight. However, it is standard in the literature to refer to many non-tight language models as language models as well. We pardon in advance the ambiguity that this introduces. Over the course of the notes, we attempt to stick to the convention that the term “language model” without qualification only refers to a tight language model whereas a “non-tight language model” is used to refer to a language model in the more colloquial sense. Linguistically, tight is acting as a non-intersective adjective. Just as in English, where a fake gun is not a gun, so too in our course notes a non-tight language model is not a language model. This distinction does in fact matter. On one hand, we can prove that many language models whose parameters are estimated from data (e.g., a finite-state language model estimated by means of maximum-likelihood estimation) are, in fact, tight. On the other hand, we can show that this is *not* true in general, i.e., *not* all language models estimated from data will be tight. For instance, a recurrent neural network language model estimated through gradient descent may not be tight (Chen et al., 2018).

When specifying p_{LM} , we have two fundamental options. Depending on whether we model $p_{\text{LM}}(\mathbf{y})$ for each string \mathbf{y} *directly* or we model *individual* conditional probabilities $p_{\text{LM}}(y_n \mid \mathbf{y}_{<n})$ we distinguish *globally* and *locally* normalized models. The names naturally come from the way the distributions in the two families are normalized: whereas globally normalized models are normalized by summing over the entire (infinite) space of strings, locally normalized models define a sequence of *conditional distributions* and make use of the chain rule of probability to define the joint probability of a whole string.

The bos symbol. Conventionally, we will include a special symbol over which globally or locally normalized models operate: the **beginning of sequence** (BOS) symbol, which, as the name suggests, denotes the beginning of a string or a sequence. For a string $\mathbf{y} = y_1 \cdots y_T$, we will suggestively denote $y_0 \stackrel{\text{def}}{=} \text{BOS}$.

2.4.1 Globally Normalized Language Models

We start with globally normalized models. Such models are also called **energy-based** language models in the literature (Bakhtin et al., 2021). To define a globally normalized language model, we start with the definition of an energy function.

Definition 2.4.1. An *energy function* is a function $\hat{p} : \Sigma^* \rightarrow \mathbb{R}$.

Inspired by concepts from statistical mechanics, an energy function can be used to define a very general class of probability distributions by normalizing its exponentiated negative values.

Now, we can define a globally normalized language model in terms of an energy function over Σ^* .

Definition 2.4.2. Let $\hat{p}_{\text{GN}}(\mathbf{y}) : \Sigma^* \rightarrow \mathbb{R}$ be an energy function. A *globally normalized model* (GNM) is defined as

$$p_{\text{LM}}(\mathbf{y}) \stackrel{\text{def}}{=} \frac{\exp[-\hat{p}_{\text{GN}}(\mathbf{y})]}{\sum_{\mathbf{y}' \in \Sigma^*} \exp[-\hat{p}_{\text{GN}}(\mathbf{y}')] } \stackrel{\text{def}}{=} \frac{1}{Z_G} \exp[-\hat{p}_{\text{GN}}(\mathbf{y})], \quad (2.10)$$

where $Z_G \stackrel{\text{def}}{=} \sum_{\mathbf{y}' \in \Sigma^*} \exp[-\hat{p}_{\text{GN}}(\mathbf{y}')]^{14}$. We call Z_G the **normalization constant**.

¹⁴We will later return to this sort of normalization when we define the softmax function in §3.1.

Globally normalized models are attractive because one only needs to define an (unnormalized) energy function \hat{p}_{GN} , which scores entire sequences at once. This is often easier than specifying a probability distribution. Furthermore, they define a probability distribution over strings $\mathbf{y} \in \Sigma^*$ *directly*. As we will see in §2.4.2, this stands in contrast to locally normalized language models which require care with the space over which they operate. However, the downside is that it may be difficult to compute the normalizer Z_G .

Normalizability

In defining the normalizer $Z_G \stackrel{\text{def}}{=} \sum_{\mathbf{y}' \in \Sigma^*} \exp[-\hat{p}_{\text{GN}}(\mathbf{y}')]$, we notationally cover up a certain subtlety. The set Σ^* is countably infinite, so Z_G may diverge to ∞ . In this case, Eq. (2.10) is not well-defined. This motivates the following definition.

Definition 2.4.3 (Normalizable energy function). *We say that an energy function is **normalizable** if the quantity Z_G in Eq. (2.10) is finite, i.e., if $Z_G < \infty$.*

With this definition, we can state a relatively trivial result that characterizes when an energy function can be turned into a globally normalized language model.

Theorem 2.4.1. *Any normalizable energy function p_{GN} induces a language model, i.e., a distribution over Σ^* .*

Proof. Given an energy function \hat{p}_{GN} , we have $\exp[-\hat{p}_{\text{GN}}(\mathbf{y})] \geq 0$ and

$$\sum_{\mathbf{y} \in \Sigma^*} p_{\text{GN}}(\mathbf{y}) = \sum_{\mathbf{y} \in \Sigma^*} \frac{\exp[-\hat{p}_{\text{GN}}(\mathbf{y})]}{\sum_{\mathbf{y}' \in \Sigma^*} \exp[-\hat{p}_{\text{GN}}(\mathbf{y}')] } \quad (2.11)$$

$$= \frac{1}{\sum_{\mathbf{y}' \in \Sigma^*} \exp[-\hat{p}_{\text{GN}}(\mathbf{y}')] } \sum_{\mathbf{y} \in \Sigma^*} \exp[-\hat{p}_{\text{GN}}(\mathbf{y})] \quad (2.12)$$

$$= 1, \quad (2.13)$$

which means that p_{GN} is a valid probability distribution over Σ^* . ■

While the fact that normalizable energy functions always form a language model is a big advantage, we will see later that *ensuring* that they are normalizable can be difficult and restrictive. This brings us to the first fundamental question of the section:

Question 2.4.1. *When is an energy function normalizable? More precisely, for which energy functions \hat{p}_{GN} is $Z_G < \infty$?*

We will not discuss any specific results here, as there are no general necessary or sufficient conditions—the answer to this of course depends on the precise definition of \hat{p}_{GN} . Later in the course notes, we will present two formalisms where we can exactly characterize when an energy function is normalizable. First, when it is weighted finite-state automaton (cf. §4.1), and, second, when it is defined through weighted context-free grammars (§4.2) and discuss the specific sufficient and necessary conditions there. However, under certain assumptions, determining whether an energy function is normalizable in the general case is undecidable.

Moreover, even if it is known that an energy function is normalizable, we still need an efficient algorithm to compute it. But, efficiently computing Z_G can be challenging: the fact that Σ^* is *infinite* means that we cannot always compute Z_G in a *tractable* way. In fact, there are no general-purpose algorithms for this. Moreover, sampling from the model is similarly intractable, as entire sequences have to be drawn at a time from the large space Σ^* .

2.4.2 Locally Normalized Language Models

The inherent difficulty in computing the normalizer, an infinite summation over Σ^* , motivates the definition of locally normalized language models, which we will denote with p_{LN} . Rather than defining a probability distribution over Σ^* directly, they decompose the problem into the problem of modeling a series of conditional distributions over the next possible symbol in the string given the context so far, i.e., $p_{\text{LN}}(y \mid \mathbf{y})$, which could

be naïvely combined into the full probability of the string by multiplying the conditional probabilities.¹⁵ Intuitively, this reduces the problem of having to normalize the distribution over an infinite set Σ^* to the problem of modeling the distribution of the *next possible symbol* y_n given the symbols seen so far $\mathbf{y}_{<n}$. This means that normalization would only ever require summation over $|\Sigma|$ symbols at a time, solving the tractability issues encountered by globally normalized models.

However, we immediately encounter another problem: In order to be a language model, $p_{LN}(\mathbf{y} \mid \mathbf{y})$ must constitute a probability distribution over Σ^* . However, as we will discuss in the next section, this may not be the case because locally normalized models can place positive probability mass on *infinitely long* sequences (cf. Example 2.5.1 in §2.5.1). Additionally, we also have to introduce a new symbol that tells us to “stop” generating a string, which we call the **end of sequence** symbol, EOS. Throughout the notes, we will assume $\text{EOS} \notin \Sigma$ and we define

$$\bar{\Sigma} \stackrel{\text{def}}{=} \Sigma \cup \{\text{EOS}\}. \quad (2.14)$$

Moreover, we will explicitly denote elements of $\bar{\Sigma}^*$ as $\bar{\mathbf{y}}$ and symbols in $\bar{\Sigma}$ as \bar{y} . Given a sequence of symbols and the EOS symbol, we take the string to be the sequence of symbols encountered *before* the *first* EOS symbol. Informally, you can think of the BOS symbol as marking the beginning of the string, and the EOS symbol as denoting the end of the string or even as a language model terminating its generation, as we will see later.

Due to the issues with defining valid probability distributions over Σ^* , we will use the term *sequence model* to refer to any model that may place positive probability on infinitely long sequences. Thus, sequence models are strictly more general than language models, which, by definition, only place positive probability mass on strings, i.e., final sequences.

Definition 2.4.4. Let Σ be an alphabet. A **sequence model** (SM) over Σ is defined as a set of conditional probability distributions

$$p_{SM}(\mathbf{y} \mid \mathbf{y}) \quad (2.15)$$

for $\mathbf{y} \in \Sigma$ and $\mathbf{y} \in \Sigma^*$. We will refer to the string \mathbf{y} in $p_{SM}(\mathbf{y} \mid \mathbf{y})$ as the **history** or the **context**.

Note that we will mostly consider SMs over the set $\bar{\Sigma}$. To reiterate, we have just formally defined locally normalized *sequence* models rather than locally normalized *language* models. That has to do with the fact that, in contrast to a globally normalized model with a normalizable energy function, a SM might not correspond to a *language* model, as alluded to at the beginning of this section and as we discuss in more detail shortly.

We will now work up to a locally normalized *language* model.

Definition 2.4.5 (Locally normalized language model). Let Σ be an alphabet. Next, let p_{SM} be a sequence model over $\bar{\Sigma}$. A **locally normalized language model** (LNM) over Σ is defined as

$$p_{LN}(\mathbf{y}) \stackrel{\text{def}}{=} p_{SM}(\text{EOS} \mid \mathbf{y}) \prod_{t=1}^T p_{SM}(y_t \mid \mathbf{y}_{<t}) \quad (2.16)$$

for $\mathbf{y} \in \Sigma^*$ with $|\mathbf{y}| = T$. We say a locally normalized language model is **tight** if

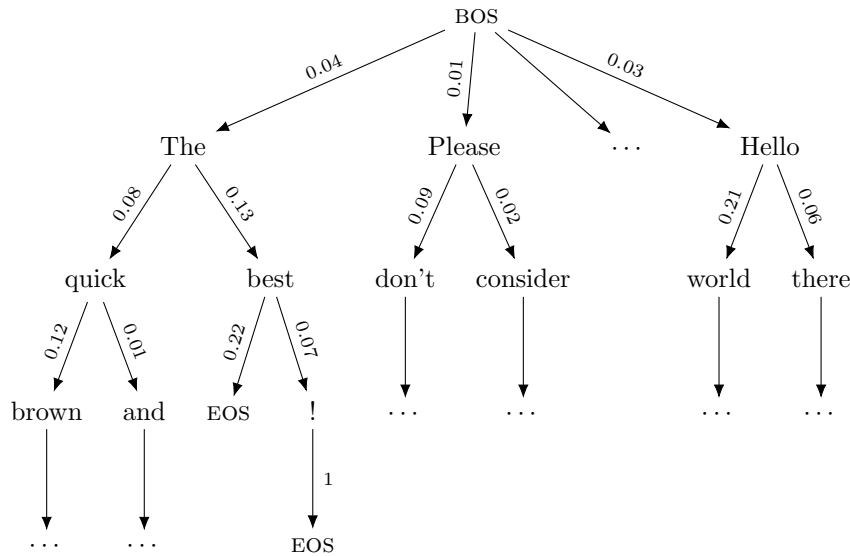
$$\sum_{\mathbf{y} \in \Sigma^*} p_{LN}(\mathbf{y}) = 1. \quad (2.17)$$

Tightness is a nuanced concept that will be discussed in great detail in §2.5.

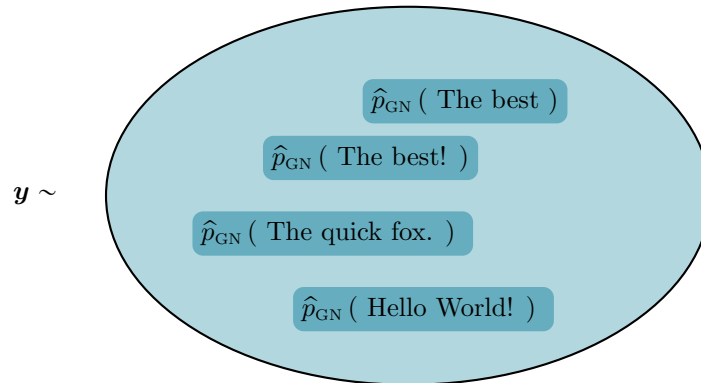
We now contrast globally and locally normalized models pictorially in the following example.

Example 2.4.1. Fig. 2.2a shows a simple instance of what a locally normalized language model would look like. We can compute the probabilities of various strings by starting at the root node BOS and choosing one of the paths to a leaf node, which will always be EOS. The values on the edges represent the conditional probabilities of observing the new word given at the target of the edge given the context seen on the path so far, i.e., $p_{LN}(\bar{y}_t \mid \bar{\mathbf{y}}_{<t})$ at the level t of the tree. For example, the probability of the string BOS “The best” EOS under this language model is $0.04 \cdot 0.13 \cdot 0.22 = 0.001144$. On the other hand, a globally normalized model would simply score all possible sentences using the score function $\hat{p}_{GN}(\mathbf{y})$, as is hinted at in Fig. 2.2b.

¹⁵We will soon see why this would not work and why we have to be a bit more careful.



(a) An example of a locally normalized language model. The values of the edges represent the conditional probability of observing the new word given the observed words (higher up on the path from the root node BOS). Note that the probabilities stemming from any inner node should sum to 1—however, to avoid clutter, only a subset of the possible arcs is drawn.



(b) An example of a globally normalized model which can for example generate sentences based on the probabilities determined by normalizing the assigned scores \hat{p}_{GN} .

Figure 2.2: “Examples” of a locally and a globally normalized language model.

Locally Normalizing a Language Model

The second fundamental question of this section concerns the relationship between language models and local normalization.

Question 2.4.2. *When can a language model be locally normalized?*

The answer to that is simple: *every* language model can be locally normalized! While the intuition behind this is very simple, the precise formulation is not. Before we discuss the details, we have to introduce the concept of prefix probabilities, which denote the sum of the probabilities of all strings beginning with a certain prefix.

Definition 2.4.6. *Let p_{LM} be a language model. We define p_{LM} 's **prefix probability** π as*

$$\pi(\mathbf{y}) \stackrel{\text{def}}{=} \sum_{\mathbf{y}' \in \Sigma^*} p_{LM}(\mathbf{y}\mathbf{y}'), \quad (2.18)$$

that is, the probability that \mathbf{y} is a prefix of any string $\mathbf{y}\mathbf{y}'$ in the language, or, equivalently, the cumulative probability of all strings beginning with \mathbf{y} .

Note that, naturally, $\pi(\varepsilon) = 1$.

Theorem 2.4.2. *Let p_{LM} be a language model. Then, there exists a locally normalized language model p_{LN} such that, for all $\mathbf{y} \in \Sigma^*$ with $|\mathbf{y}| = T$,*

$$p_{LM}(\mathbf{y}) = p_{LN}(\mathbf{y}) = p_{SM}(\text{EOS} \mid \mathbf{y}) \prod_{t=1}^T p_{SM}(y_t \mid \mathbf{y}_{<t}). \quad (2.19)$$

Proof. We define the individual conditional probability distributions over the next symbol of the SM p_{SM} using the chain rule of probability. If $\pi(\mathbf{y}) > 0$, then define

$$p_{SM}(y \mid \mathbf{y}) \stackrel{\text{def}}{=} \frac{\pi(\mathbf{y}y)}{\pi(\mathbf{y})} \quad (2.20)$$

for $y \in \Sigma$ and $\mathbf{y} \in \Sigma^*$ such that $p(\mathbf{y}) > 0$. We still have to define the probabilities of *ending* the sequence using p_{SM} by defining the EOS probabilities. We define, for any $\mathbf{y} \in \Sigma^*$ such that $\pi(\mathbf{y}) > 0$,

$$p_{SM}(\text{EOS} \mid \mathbf{y}) \stackrel{\text{def}}{=} \frac{p_{LM}(\mathbf{y})}{\pi(\mathbf{y})} \quad (2.21)$$

that is, the probability that the globally normalized model will generate *exactly* the string \mathbf{y} and not any continuation of it $\mathbf{y}\mathbf{y}'$, given that \mathbf{y} has already been generated. Each of the conditional distributions of this model (Eqs. (2.20) and (2.21)) is clearly defined over $\bar{\Sigma}$. This, therefore, defines a valid SM. To see that p_{LN} constitutes the same distribution as p_{LM} , consider two cases.

Case 1: Assume $\pi(\mathbf{y}) > 0$. Then, we have

$$p_{LN}(\mathbf{y}) = \left[\prod_{t=1}^T p_{SM}(y_t \mid \mathbf{y}_{<t}) \right] p_{SM}(\text{EOS} \mid \mathbf{y}) \quad (2.22)$$

$$\begin{aligned} &= \frac{\pi(\cancel{\mathbf{y}_1})}{\pi(\varepsilon)} \frac{\pi(\cancel{\mathbf{y}_1}\cancel{\mathbf{y}_2})}{\pi(\cancel{\mathbf{y}_1})} \dots \frac{\pi(\cancel{\mathbf{y}_{<T}})}{\pi(\cancel{\mathbf{y}_{<T-1}})} \frac{\pi(\mathbf{y})}{\pi(\cancel{\mathbf{y}_{<T}})} p_{SM}(\text{EOS} \mid \mathbf{y}) \\ &= \frac{\pi(\cancel{\mathbf{y}})}{\pi(\varepsilon)} \frac{p_{LM}(\mathbf{y})}{\pi(\cancel{\mathbf{y}})} \end{aligned} \quad (2.23)$$

$$= p_{LM}(\mathbf{y}) \quad (2.24)$$

where $\pi(\varepsilon) = 1$.

Case 2: Assume $\pi(\mathbf{y}) = 0$. Let $\mathbf{y} = y_1 \cdots y_T$. Then, there must exist a $1 \leq t' \leq T$ such that $\pi(\mathbf{y}_{<t'}) = 0$. Note that

$$p_{\text{LN}}(\mathbf{y}) = \prod_{t=1}^{t'} p_{\text{SM}}(y_t \mid \mathbf{y}_{<t}) = 0 \quad (2.25)$$

whereas the conditional probabilities after t' can be arbitrarily defined since they do not affect the string having 0 probability. ■

When Is a Locally Normalized Language Model a Language Model?

LNMs which specify distributions over strings $p_{\text{LN}}(y_1 \dots y_T)$ in terms of their conditional probabilities $p_{\text{SM}}(y_t \mid \mathbf{y}_{<t})$ for $t = 1, \dots, T$ and $p_{\text{SM}}(\text{EOS} \mid \mathbf{y})$ have become the standard in NLP literature. However, LNMs come with their own set of problems. An advantage of normalizable globally normalized models is that they, by definition, always define a *valid* probability space over $\bar{\Sigma}$. Although this might be counterintuitive at first, the same cannot be said for LNMs—in this sense, locally normalized “language models” might not even be language models! One might expect that in a LNM p_{LN} , it would hold that $\sum_{\mathbf{y} \in \Sigma^*} p_{\text{LN}}(\mathbf{y}) = 1$. However, this might not be the case! This is the issue with the terminology we brought up earlier and it brings us to the last fundamental question of this section.

Question 2.4.3. *When does an LNM encode a language model?*

As the conditions are a bit more nuanced, it requires a longer treatment. We explore this issue in much more detail in the next section.

2.5 Tight Language Models

We saw in the last section that any language model p_{LM} can be converted into a locally normalized sequence model (cf. §2.4.2). The converse, however, is *not* true. As alluded to in the previous section and as we detail in this section, there exist sets of conditional distributions $p_{\text{LN}}(\bar{y} \mid \bar{\mathbf{y}})$ over $\bar{\Sigma}^*$ such that $p_{\text{LN}}(\bar{\mathbf{y}})$ as defined in Eq. (2.15) does not represent a valid probability measure over Σ^* (after taking into account the semantics of EOS), i.e., over the set of *finite* strings. Indeed, we will later show that some popular classes of locally normalized sequence models used in practice have parameter settings in which the generative process terminates with probability < 1 . This means that p_{LN} “leaks” some of its probability mass to *infinite* sequences. This section investigates this behavior in a lot of detail. It is based on the recent work from Du et al. (2022).

2.5.1 Tightness

Models whose generative process may fail to terminate are called **non-tight** (Chi, 1999).¹⁶

Definition 2.5.1. A locally normalized language model p_{LN} derived from a sequence model p_{SM} is called **tight** if it defines a valid probability distribution over Σ^* :

$$\sum_{\mathbf{y} \in \Sigma^*} p_{\text{LN}}(\mathbf{y}) = \sum_{\mathbf{y} \in \Sigma^*} \left[p_{\text{SM}}(\text{EOS} \mid \mathbf{y}) \prod_{t=1}^T p_{\text{SM}}(y_t \mid \mathbf{y}_{<t}) \right] = 1. \quad (2.26)$$

Note that the individual conditional distributions $p_{\text{SM}}(y \mid \mathbf{y})$ in a non-tight LNM still are valid conditional distributions (i.e., they sum to one). However, the distribution over all possible strings that they induce may not sum to 1. To be able to investigate this phenomenon more closely, let us first examine what the conditional probabilities of an LNM actually define and how they can result in non-tightness. We now ask ourselves: given a sequence model p_{SM} , what is p_{LN} ? Is p_{LN} a language model, i.e., a distribution over Σ^* (after taking into account the semantics of EOS)? Certainly, the answer is yes if the LNM’s conditional probabilities match the conditional probabilities of some known language model p_{LM} as defined in §2.4.2,¹⁷ in which case p_{LN} is specifically the language model p_{LM} itself. In this case clearly $p_{\text{LN}}(\Sigma^*) \stackrel{\text{def}}{=} \sum_{\mathbf{y} \in \Sigma^*} p_{\text{LN}}(\mathbf{y}) = \sum_{\mathbf{y} \in \Sigma^*} p_{\text{LM}}(\mathbf{y}) = 1$. If instead $p_{\text{LN}}(\Sigma^*) < 1$, the LNM’s conditional probabilities do *not* match the conditional probabilities of any language model p_{LM} .

To see how this can happen, we now exhibit such an LNM in the following example.

Example 2.5.1. Consider the bigram model defined in Fig. 2.3a over the alphabet $\Sigma = \{a, b\}$.¹⁸ Although the conditional probability distributions $p_{\text{LN}}(\cdot \mid \mathbf{y}_{<n})$ each sum to 1 over $\bar{\Sigma}$, they fail to combine into a model p_{LN} that sums to 1 over Σ^* (i.e., a language model): under this model, any finite string that contains the symbol b will have probability 0, since $p_{\text{LN}}(\text{EOS} \mid b) = p_{\text{LN}}(a \mid b) = 0$. This implies $p_{\text{LN}}(\Sigma^*) = \sum_{n=0}^{\infty} p_{\text{LN}}(a^n) = \sum_{n=0}^{\infty} (0.7)^n \cdot 0.1 = \frac{0.1}{1-0.7} = \frac{1}{3} < 1$.

Example 2.5.2. On the other hand, in the bigram model in Fig. 2.3b, obtained from Example 2.5.1 by

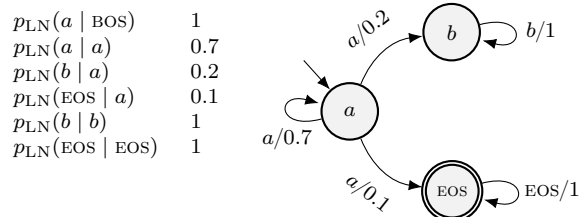
¹⁶Tight models are also called **consistent** (Booth and Thompson, 1973; Chen et al., 2018) and **proper** (Chi, 1999) in the literature.

¹⁷That is, $p_{\text{LM}}(y_t \mid \mathbf{y}_{<t}) = p_{\text{LN}}(y_t \mid \mathbf{y}_{<t})$ whenever the former conditional probability is well-defined under the language model p_{LM} , i.e., whenever $y_t \in \bar{\Sigma}$ and $\mathbf{y}_{<t} \in \Sigma^*$ with $p_{\text{LM}}(\mathbf{y}_{<t}) > 0$.

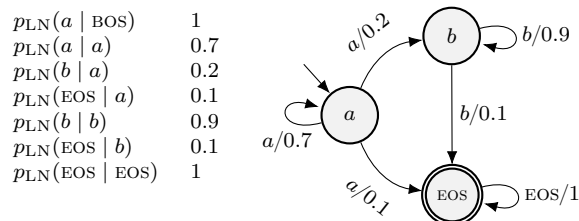
¹⁸The graphical representation of the LNM depicts a so-called weighted finite-state automaton, a framework of language models we will introduce shortly. For now, it is not crucial that you understand the graphical representation and you can simply focus on the conditional probabilities specified in the figure.

changing the arcs from the b state, $p_{LN}(\Sigma^*) = 1$. We can see that by calculating:

$$\begin{aligned}
 \mathbb{P}(\Sigma^*) &= \sum_{n=1}^{\infty} \sum_{m=0}^{\infty} \mathbb{P}(a^n b^m) \\
 &= \sum_{n=1}^{\infty} \left(\mathbb{P}(a^n) + \sum_{m=1}^{\infty} \mathbb{P}(a^n b^m) \right) \\
 &= \sum_{n=1}^{\infty} \left(0.1 \cdot (0.7)^{n-1} + \sum_{m=1}^{\infty} (0.7)^{n-1} \cdot 0.2 \cdot (0.9)^{m-1} \cdot 0.1 \right) \\
 &= \sum_{n=1}^{\infty} \left(0.1 \cdot (0.7)^{n-1} + (0.7)^{n-1} \cdot 0.2 \cdot \frac{1}{1-0.9} \cdot 0.1 \right) \\
 &= \sum_{n=1}^{\infty} (0.1 \cdot (0.7)^{n-1} + 0.2 \cdot (0.7)^{n-1}) \\
 &= \sum_{n=1}^{\infty} 0.3 \cdot (0.7)^{n-1} = \frac{0.3}{1-0.7} = 1.
 \end{aligned}$$



(a) Non-tight 2-gram model.



(b) Tight 2-gram model.

Figure 2.3: Tight and non-tight bigram models, expressed as Mealy machines. Symbols with conditional probability of 0 are omitted.

Example 2.5.1 confirms that the local normalization does not necessarily yield p_{LN} that is a valid distribution over Σ^* . But if p_{LN} is not a language model, *what* is it? It is intuitive to suspect that, in a model with $p_{LN}(\Sigma^*) < 1$, the remainder of the probability mass “leaks” to infinite sequences, i.e., the generative process may continue forever with probability > 0 . This means that, to be able to characterize p_{LN} , we will have to be able to somehow take into account infinite sequences. We will make this intuition formal below.

Delving a bit deeper, the non-tightness of Example 2.5.1 is related to the fact that the conditional probability of EOS is 0 at some states, in contrast to Example 2.5.2. However, requiring $p_{LN}(y_n = \text{EOS} \mid \mathbf{y}_{<n}) > 0$ for all prefixes $\mathbf{y}_{<n}$ is neither *necessary* nor *sufficient* to ensure tightness. It is not necessary because one can, for example, construct an LNM in which $p_{LN}(y_n = \text{EOS} \mid \mathbf{y}_{<n}) = 0.1$ when n is even but $= 0$ otherwise. Such a model generates only odd-length strings but is tight. We will postpone non-sufficiency for later, where we will present specific LNMs under which the conditional probability of EOS is always > 0 , yet are non-tight.

2.5.2 Defining the probability measure of an LNM

We now rigorously characterize the kind of distribution induced by an LNM, i.e., we investigate what p_{LNM} is. As mentioned earlier, an LNM can lose probability mass to the set of infinite sequences, Σ^∞ . However, Σ^∞ , unlike Σ^* , is *uncountable*, and it is due to this fact that we need to work explicitly with the *measure-theoretic* formulation of probability which we introduced in §2.2. We already saw the peril of not treating distributions over uncountable sets carefully is necessary in Example 2.1.1—the set of all infinite sequences of coin tosses is indeed uncountable.

Including infinite strings and eos. As we saw in Example 2.1.1, sampling successive symbols from a non-tight LNM has probability > 0 of continuing forever, i.e., generating infinite strings. Motivated by that, we hope to regard the LNM as defining a valid probability space over $\Omega = \Sigma^* \cup \Sigma^\infty$, i.e., both finite as well as infinite strings, and then “relate” it to our definition of true language models. Notice, however, that we also have to account for the difference in the alphabets: while we would like to characterize language models in terms of strings over the alphabet Σ , LNMs work over symbols in $\bar{\Sigma}$.

With this in mind, we now embark on our journey of discovering what p_{LNM} represents. Given an LNM, we will first need to turn its p_{LNM} into a measurable space by defining an appropriate σ -algebra. This type of distribution is more general than a language model as it works over both finite as well as infinite sequences. To distinguish the two, we will expand our vocabulary and explicitly *differentiate* between true language models and non-tight LNMs. We will refer to a distribution over $\Sigma^* \cup \Sigma^\infty$ as a sequence model. As noted in our definition of a sequence model (cf. Definition 2.4.4), an LNM defines a probability measure over $\Sigma^* \cup \Sigma^\infty$. Thus, an equivalent distribution, which will be useful for this section, would be the following.

Definition 2.5.2 (Sequence model). *A **sequence model** is a probability space over the set $\Sigma^* \cup \Sigma^\infty$.*

Intuitively, and we will make this precise later, the set $\Sigma^\infty \subset \Sigma^* \cup \Sigma^\infty$ in Definition 2.5.2 represents the *event* where the sequence model is **non-terminating**, i.e., it attempts to generate an infinitely long sequence. We can then understand language models in a new sense.

Definition 2.5.3 (Re-definition of a Language model). *A **language model** is a probability space over Σ^* . Equivalently, a language model is a sequence model such that $\mathbb{P}(\Sigma^\infty) = 0$.*

Now buckle up! Our goal through the rest of this section is to rigorously construct a probability space of a sequence model as in Definition 2.2.2 and Definition 2.5.2 which encodes the probabilities assigned by an LNM. Then, we will use this characterization to formally investigate tightness. An outline of what this is going to look like is shown in Fig. 2.4.

Defining an Algebra over $\bar{\Sigma}^\infty$ (Step 1)

Since an LNM produces conditional distributions over the augmented alphabet $\bar{\Sigma}$ (first box in Fig. 2.4) and results in possibly infinite strings, we will first construct a probability space over $\bar{\Sigma}^\infty$, which will naturally induce a sequence model. We will do that by first constructing an *algebra* (cf. Definition 2.2.3) over $\Omega = \bar{\Sigma}^\infty$ for some alphabet Σ (second box in Fig. 2.4). Then, assuming we are given an LNM p_{LNM} over $\bar{\Sigma}$, we will associate the constructed algebra with a pre-measure (cf. Definition 2.2.4) that is “consistent” with p_{LNM} (third box in Fig. 2.4).

We will make use of the following definition to construct the algebra:

Definition 2.5.4. *Given any set $\mathcal{H} \subseteq \bar{\Sigma}^k$, i.e., a set of sequences of symbols from $\bar{\Sigma}$ of length k , define its **cylinder set** (of rank k) to be*

$$\bar{C}(\mathcal{H}) \stackrel{\text{def}}{=} \{\mathbf{y}\omega : \mathbf{y} \in \mathcal{H}, \omega \in \bar{\Sigma}^\infty\} \quad (2.27)$$

In essence, a cylinder set of rank k is the set of infinite strings that share their k -prefix with some string $\bar{\mathbf{y}} \in \mathcal{H} \subseteq \bar{\Sigma}^k$. In particular, for a length- k string $\bar{\mathbf{y}} = \bar{y}_1 \cdots \bar{y}_k$, the cylinder set $\bar{C}(\bar{\mathbf{y}}) \stackrel{\text{def}}{=} \bar{C}(\{\bar{\mathbf{y}}\})$ is the set of all infinite strings prefixed by $\bar{\mathbf{y}}$.¹⁹

¹⁹This type of cylinder set, i.e., one that is generated by a singleton, is also called a **thin cylinder**.

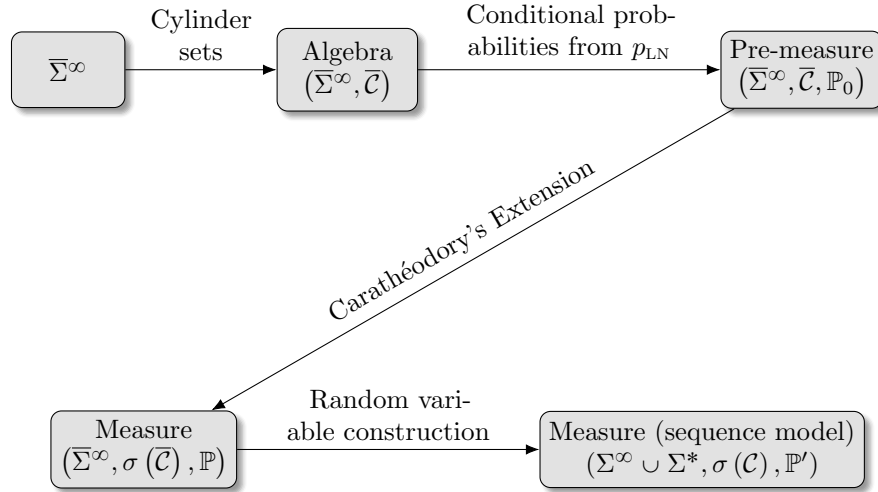


Figure 2.4: The outline of our measure-theoretic treatment of LNM in this section to arrive at a precise characterization of p_{LN} . The final box corresponds to the sequence model (probability measure over $\Sigma^* \cup \Sigma^\infty$) constructed for p_{LN} .

We denote the collection of all rank- k cylinder sets by

$$\bar{\mathcal{C}}_k \stackrel{\text{def}}{=} \{\bar{\mathcal{C}}(\mathcal{H}) : \mathcal{H} \in \mathcal{P}(\bar{\Sigma}^k)\} \quad (2.28)$$

and define

$$\bar{\mathcal{C}} \stackrel{\text{def}}{=} \bigcup_{k=1}^{\infty} \bar{\mathcal{C}}_k \quad (2.29)$$

to be the collection of all cylinder sets over Ω .²⁰

The following lemma asserts $\bar{\mathcal{C}} \subseteq \mathcal{P}(\Omega)$ is what we want in the second block of Fig. 2.4.

Lemma 2.5.1. $\bar{\mathcal{C}} \subseteq \mathcal{P}(\Omega)$ is an algebra over $\Omega = \bar{\Sigma}^\infty$.

Proof. First, $\bar{\Sigma}^\infty = \bar{\mathcal{C}}(\bar{\Sigma}^k)$ for any k , and in particular is a cylinder set of any rank. Secondly, given a cylinder set $\bar{\mathcal{C}}(\mathcal{H})$ of rank k , i.e., $\mathcal{H} \subseteq \bar{\Sigma}^k$, $(\bar{\mathcal{C}}(\mathcal{H}))^c = \bar{\mathcal{C}}(\bar{\Sigma}^k \setminus \mathcal{H})$. Hence, $\bar{\mathcal{C}}$ is closed under complements. Finally, notice that the intersection of two cylinder sets of ranks $k_1 \leq k_2$ is another cylinder set of rank k_2 . Hence, $\bar{\mathcal{C}}$ is an algebra over Ω . ■

With this, the first step of Fig. 2.4 is done!

Defining a Pre-measure over $\bar{\mathcal{C}}$ (Step 2)

We are now ready to define the pre-measure \mathbb{P}_0 for the cylinder algebra $\bar{\mathcal{C}}$. Given an LNM p_{LN} and any set $\bar{\mathcal{C}}(\mathcal{H}) \in \bar{\mathcal{C}}$, let

$$\mathbb{P}_0(\bar{\mathcal{C}}(\mathcal{H})) \stackrel{\text{def}}{=} \sum_{\bar{\mathbf{y}} \in \mathcal{H}} p_{\text{LN}}(\bar{\mathbf{y}}) \quad (2.30)$$

where we have defined

$$p_{\text{LN}}(\bar{\mathbf{y}}) \stackrel{\text{def}}{=} \prod_{t=1}^T p_{\text{LN}}(\bar{y}_t \mid \bar{\mathbf{y}}_{<t}). \quad (2.31)$$

Note that there is a caveat here since the same cylinder set may admit different \mathcal{H} .²¹ Before showing that \mathbb{P}_0 defines a valid pre-measure, we address this and show that \mathbb{P}_0 is indeed well defined.

²⁰We invite the reader to verify that $\bar{\mathcal{C}}_1 \subset \bar{\mathcal{C}}_2 \subset \bar{\mathcal{C}}_3 \subset \dots$.

²¹For example, in the infinite coin toss model, $\mathcal{C}(\mathbf{H}) = \mathcal{C}(\{\mathbf{HH}, \mathbf{HT}\})$.

Proposition 2.5.1. \mathbb{P}_0 as defined in Eq. (2.30) is a well-defined function.

Proof. Suppose a cylinder set can be described by two different prefix sets: $H_1 \subseteq \bar{\Sigma}^{k_1}$ and $H_2 \subseteq \bar{\Sigma}^{k_2}$. In other words, $\bar{C}(H_1) = \bar{C}(H_2)$. Without loss of generality, assume that $k_1 \leq k_2$. Then,

$$\bar{C}(H_2) = \bar{C}(H_1) \tag{2.32a}$$

$$= \bigcup_{\mathbf{y} \in H_1} \bar{C}(\mathbf{y}) \tag{2.32b}$$

$$= \bigcup_{\mathbf{y} \in H_1} \bigcup_{\bar{\mathbf{y}} \in \bar{\Sigma}^{k_2-k_1}} \bar{C}(\mathbf{y}\bar{\mathbf{y}}). \tag{2.32c}$$

All the unions above are disjoint, and hence $H_2 = \bigcup_{\mathbf{y} \in \bar{\Sigma}^{k_2-k_1}} \{\mathbf{y}\mathbf{y} : \mathbf{y} \in H_1\}$. Then, by the locally-normalizing property of p_{LN} , we have that

$$\mathbb{P}_0(\bar{C}(H_1)) = \mathbb{P}_0(\bar{C}(H_2)). \tag{2.33}$$

■

With this, we are able to state and prove the lemma which shows that \mathbb{P}_0 is a pre-measure, which is what we need in the third block of Fig. 2.4.

Lemma 2.5.2. \mathbb{P}_0 is a pre-measure over $\bar{\mathcal{C}}$.

For the proof of Lemma 2.5.2, we will mostly follow the proof of Theorem 2.3 in Billingsley (1995), with the exception of invoking the Tychonoff theorem directly. This proof depends on the following lemma, which is Example 2.10 in Billingsley (1995). We repeat the statement and proof here for the reader's convenience.

Lemma 2.5.3. Let \mathbb{P}_0 be a finitely additive probability pre-measure over $\bar{\mathcal{C}}$ such that, given a decreasing sequence of sets $A_1 \supset A_2 \supset \dots$ in $\bar{\mathcal{C}}$ where $\bigcap_{n=1}^{\infty} A_n = \emptyset$, $\lim_{n \rightarrow \infty} \mathbb{P}_0(A_n) = 0$. Then, \mathbb{P}_0 is also countably additive over $\bar{\mathcal{C}}$.

Proof. Let $\{A_n\}$ be a sequence of disjoint sets in $\bar{\mathcal{C}}$ such that $A = \bigcup_n A_n \in \bar{\mathcal{C}}$. Then, defining $B_n = \bigcup_{m>n} A_m$, we see that $B_1 \supset B_2 \supset \dots$ and $\bigcap_n B_n = \emptyset$. Notice that

$$A = A_1 \cup B_1 = A_1 \cup A_2 \cup B_2 = \dots = A_1 \cup \dots \cup A_n \cup B_n \tag{2.34}$$

for any n and hence by finite additivity of \mathbb{P}_0

$$\mathbb{P}_0(A) = \mathbb{P}_0(A_1) + \dots + \mathbb{P}_0(A_n) + \mathbb{P}_0(B_n) \tag{2.35}$$

or equivalently

$$\mathbb{P}_0(A_1) + \dots + \mathbb{P}_0(A_n) = \mathbb{P}_0(A) - \mathbb{P}_0(B_n). \tag{2.36}$$

Since, $B_n \downarrow \emptyset$ implies that $\mathbb{P}_0(B_n) \downarrow 0$ by assumption, taking the limits on both sides of Eq. (2.36) yields

$$\sum_n \mathbb{P}_0(A_n) = \lim_{n \rightarrow \infty} \sum_{i \leq n} \mathbb{P}_0(A_i) = \mathbb{P}_0(A) - \lim_{n \rightarrow \infty} \mathbb{P}_0(B_n) = \mathbb{P}_0(A) \tag{2.37}$$

which shows countable additivity. ■

We also recall the Tychonoff theorem.²²

Theorem 2.5.1 (Tychonoff). Let $\{\mathcal{X}_\alpha\}_{\alpha \in J}$ be an indexed family of compact topologies. Then, their product topology $\prod_{\alpha \in J} \mathcal{X}_\alpha$ is also compact.

We can now give the proof for Lemma 2.5.2.

²²See §37 in Munkres (2000) for a detailed and well-written treatise.

Proof of Lemma 2.5.2. We first show that \mathbb{P}_0 is finitely additive over $\bar{\mathcal{C}}$. Let $\mathcal{C}(\mathcal{H}_1)$ and $\mathcal{C}(\mathcal{H}_2)$ be two disjoint cylinder sets. By Proposition 2.5.1, we can assume they are of the same rank without loss of generality. Then,

$$\mathcal{C}(\mathcal{H}_1) \cup \mathcal{C}(\mathcal{H}_2) = \bigcup_{\mathbf{y} \in \mathcal{H}_1} \{\mathbf{y}\omega : \omega \in \bar{\Sigma}^\infty\} \cup \bigcup_{\mathbf{y} \in \mathcal{H}_2} \{\mathbf{y}\omega : \omega \in \bar{\Sigma}^\infty\} \quad (2.38a)$$

$$= \bigcup_{\mathbf{y} \in \mathcal{H}_1 \cup \mathcal{H}_2} \{\mathbf{y}\omega : \omega \in \bar{\Sigma}^\infty\} \quad (\mathcal{H}_1 \text{ and } \mathcal{H}_2 \text{ equal rank and disjoint}) \quad (2.38b)$$

$$= \mathcal{C}(\mathcal{H}_1 \cup \mathcal{H}_2) \quad (2.38c)$$

which leads to

$$\mathbb{P}_0(\mathcal{C}(\mathcal{H}_1) \cup \mathcal{C}(\mathcal{H}_2)) = \mathbb{P}_0(\mathcal{C}(\mathcal{H}_1 \cup \mathcal{H}_2)) \quad (2.39a)$$

$$= \sum_{\mathbf{y} \in \mathcal{H}_1 \cup \mathcal{H}_2} p_{\text{LN}}(\mathbf{y}) \quad (2.39b)$$

$$= \mathbb{P}_0(\mathcal{C}(\mathcal{H}_1)) + \mathbb{P}_0(\mathcal{C}(\mathcal{H}_2)). \quad (2.39c)$$

Hence, \mathbb{P}_0 is finitely additive.

Now, equip $\bar{\Sigma}$ with the discrete topology. Since $\bar{\Sigma}$ is finite, it is compact under the discrete topology and so is $\bar{\Sigma}^\infty$ by Theorem 2.5.1. Then, by properties of the product topology over discrete finite spaces, all cylinder sets in $\bar{\Sigma}^\infty$ are compact. To apply Lemma 2.5.3, let $\mathcal{C}_1 \supset \mathcal{C}_2 \supset \dots$ be a decreasing sequence of cylinder sets with empty intersection. Suppose to the contrary that $\mathbb{P}_0(\bigcap_n \mathcal{C}_n) > 0$. This would imply that all \mathcal{C}_n are nonempty (any of these being empty would result in a measure 0). However, by Cantor's intersection theorem²³, $\bigcap_n \mathcal{C}_n$ is nonempty, contradicting the assumption. Hence, $\mathbb{P}_0(\bigcap_n \mathcal{C}_n) = 0$, and by Lemma 2.5.3, \mathbb{P}_0 is countably additive.

With this, we have proved that \mathbb{P}_0 is countably additive. To show that \mathbb{P}_0 defines a pre-measure, we still have to show that $\mathbb{P}_0(\Omega) = 1$. Recall from the proof of Lemma 2.5.1 that $\bar{\Sigma}^\infty = \bar{C}(\bar{\Sigma}^k)$ for any $k > 0$. In particular, $\bar{\Sigma}^\infty = \bar{C}(\bar{\Sigma}^1) = \bar{C}(\bar{\Sigma})$. This means that

$$\mathbb{P}_0(\Omega) = \mathbb{P}_0(\bar{C}(\bar{\Sigma})) \quad (2.40)$$

$$= \sum_{\bar{\mathbf{y}} \in \bar{\Sigma}} p_{\text{LN}}(\bar{\mathbf{y}}) \quad (2.41)$$

$$= \sum_{\bar{\mathbf{y}} \in \bar{\Sigma}} p_{\text{LN}}(\bar{\mathbf{y}} \mid \text{BOS}) = 1. \quad (2.42)$$

The last equality follows from local normalization of the sequence model. ■

With this, we have successfully completed the first two steps of Fig. 2.4! However, we have only defined a *pre-measure* over the set of *infinite* EOS-containing sequences $\bar{\Sigma}^\infty$. This does not yet satisfy all the properties we would like from a probability space. Because of that, we next *extend* the constructed probability pre-measure \mathbb{P}_0 into a valid probability measure \mathbb{P} to arrive to a valid probability space.

Extending the Pre-measure \mathbb{P}_0 into a Measure \mathbb{P} (Step 3)

To extend \mathbb{P}_0 into a measure, we will use Carathéodory's theorem:

Theorem 2.5.2 (Carathéodory's Extension Theorem). *Given an algebra \mathcal{A} over some set Ω and a probability pre-measure $\mathbb{P}_0 : \mathcal{A} \rightarrow [0, 1]$, there exists a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ such that $\mathcal{A} \subset \mathcal{F}$ and $\mathbb{P}|_{\mathcal{A}} = \mathbb{P}_0$. Furthermore, the σ -algebra \mathcal{F} depends only on \mathcal{A} and is minimal and unique, which we will also denote by $\sigma(\mathcal{A})$, and the probability measure \mathbb{P} is unique.*

²³Cantor's intersection theorem states that a decreasing sequence of nonempty compact sets have a nonempty intersection. A version of this result in introductory real analysis is the Nested Interval Theorem.

Proof Sketch. First, construct an outer measure by approximation with countable coverings. Then, show that the collection of sets that is measurable with respect to this outer measure is a σ -algebra \mathcal{F} that contains \mathcal{A} . Finally, restricting the outer measure to this σ -algebra, one is then left with a probability space. To show minimality, one can show that \mathcal{F} is contained in any σ -algebra that contains \mathcal{A} . Uniqueness is given by applying Dynkin's π - λ theorem (Theorem 3.2 in Billingsley, 1995).

Great care must be taken in each step involved in the outline above. To address these is well beyond the scope of this treatment and we refer the reader to the many excellent texts with a proof of this theorem, such as Chapter 12 in Royden (1988) and Chapter 11 in Billingsley (1995). ■

Applying Carathéodory's extension theorem to our cylinder algebra $\bar{\mathcal{C}}$ and pre-measure \mathbb{P}_0 , we see that there exists a probability space $(\bar{\Sigma}^\infty, \sigma(\bar{\mathcal{C}}), \mathbb{P})$ over $\bar{\Sigma}^\infty$ that agrees with the LNM p_{LN} 's probabilities.

Phew! This now gets us to the fourth box in Fig. 2.4 and we only have one step remaining.

Defining a Sequence Model (Step 4)

We now have to make sure that the *outcome space* of the defined probability space fits the definition of a sequence model. That is, we have to find a way to convert (map) the infinite EOS-containing sequences from $\bar{\Sigma}^\infty$ into EOS-free finite or possibly infinite strings processed by a sequence model as required by Definition 2.5.2. We will achieve this through the use of a *random variable*.

Recall from Definition 2.2.5 that a random variable is a mapping between *two* sigma-algebras. Since we want our final measure space to work with the outcome space $\Sigma^* \cup \Sigma^\infty$, we, therefore, want to construct a σ -algebra over $\Sigma^* \cup \Sigma^\infty$ and then map elements from $\bar{\Sigma}^\infty$ to $\Sigma^* \cup \Sigma^\infty$ to have the appropriate objects. We will do so in a similar fashion as we constructed $(\bar{\Sigma}^\infty, \bar{\mathcal{C}})$. Given $\mathcal{H} \subset \Sigma^k$, define a rank- k cylinder set in $\Sigma^* \cup \Sigma^\infty$ to be

$$\mathcal{C}(\mathcal{H}) \stackrel{\text{def}}{=} \{\mathbf{y}\omega : \mathbf{y} \in \mathcal{H}, \omega \in \Sigma^* \cup \Sigma^\infty\}. \quad (2.43)$$

Notice the major change from Eq. (2.27): the suffixes ω of the elements in $\mathcal{C}(\mathcal{H})$ now come from $\Sigma^* \cup \Sigma^\infty$ rather than $\bar{\Sigma}^\infty$. This means (i) that they do not contain EOS and (ii) that they (and thus, elements of $\mathcal{C}(\mathcal{H})$) can also be finite. Let \mathcal{C}_k be the set of all rank- k cylinder sets. Define $\mathcal{C} \stackrel{\text{def}}{=} \bigcup_{k=1}^\infty \mathcal{C}_k$. Then, $\sigma(\mathcal{C})$ is a σ -algebra by the same reasoning as in Lemma 2.5.1 and Theorem 2.5.2. We can now define the following random variable

$$x(\omega) = \begin{cases} \omega_{<k} & \text{if } k \text{ is the first EOS in } \omega, \\ \omega & \text{otherwise (if EOS } \notin \omega) \end{cases} \quad (2.44)$$

given any $\omega \in \bar{\Sigma}^\infty$. The proposition below shows that x is well-defined.

Proposition 2.5.2. *The function $x : (\bar{\Sigma}^\infty, \sigma(\bar{\mathcal{C}})) \rightarrow (\Sigma^* \cup \Sigma^\infty, \sigma(\mathcal{C}))$ defined in Eq. (2.44) is a measurable mapping.*

Proof. To show that x is measurable, it suffices to show the measurability of preimage of a generating set of the σ -algebra. Note that the set of thin cylinder sets is a generating set. Let $\mathcal{C}(\mathbf{y})$ be a thin cylinder set,

$$x^{-1}(\mathcal{C}(\mathbf{y})) = x^{-1}(\{\mathbf{y}\omega : \omega \in \Sigma^* \cup \Sigma^\infty\}) \quad (2.45a)$$

$$= x^{-1}(\{\mathbf{y}\omega : \omega \in \Sigma^*\}) \cup x^{-1}(\{\mathbf{y}\omega : \omega \in \Sigma^\infty\}) \quad (2.45b)$$

$$= \left(\bigcup_{\omega \in \Sigma^*} \bar{\mathcal{C}}(\mathbf{y}\omega\text{EOS}) \right) \cup \left(\bar{\mathcal{C}}(\mathbf{y}) \cap \bigcap_{k=1}^\infty \mathcal{A}_k^c \right) \quad (2.45c)$$

Note that the sets \mathcal{A}_k above are defined in Eq. (2.58) which are cylinder sets representing the event of terminating at step k . Then, from the derivation above, we can see that $x^{-1}(\mathcal{C}(\mathbf{y}))$ is formed by countable operations over measurable sets (cylinder sets) in $\bar{\Sigma}^\infty$, and is hence measurable. So x is a measurable function. ■

x intuitively “cuts out” the first stretch of ω before the first EOS symbol (where an LNM would stop generating) or leaves the sequence intact if there is no termination symbol EOS. One can check that \mathbb{P}_* ,

defined using \mathbb{P} , is indeed a probability measure on $(\Sigma^* \cup \Sigma^\infty, \sigma(\mathcal{C}))$ and hence $(\Sigma^* \cup \Sigma^\infty, \sigma(\mathcal{C}), \mathbb{P}_*)$ is a probability space. We have therefore arrived at the final box of Fig. 2.4 and shown that, given any LNM, we can construct an associated sequence model as defined in Definition 2.5.2! In other words, given an LNM p_{LN} , we have constructed a sequence model p_{SM} (a probability space over $\Sigma^\infty \cup \Sigma^*$ where the probabilities assigned to (infinite) strings by p_{SM} agree with p_{LN}).

2.5.3 Interpreting the Constructed Probability Space

Under the formulation of a probability space together with a random variable, useful probability quantities arise naturally and intuitively.

Consider, for example, the probability of a single finite string $\mathbf{y} \in \Sigma^*$, $\mathbb{P}_*(\mathbf{y})$. By definition of \mathbf{x} , this equals

$$\mathbb{P}_*(\mathbf{y}) = \mathbb{P}_*(\mathbf{x} = \mathbf{y}) \quad (2.46)$$

$$= \mathbb{P}(\mathbf{x}^{-1}(\mathbf{y})) \quad (2.47)$$

$$= \mathbb{P}(\text{All the sequences } \boldsymbol{\omega} \in \bar{\Sigma}^\infty \text{ which map to } \mathbf{y}.) \quad (2.48)$$

All the sequences $\boldsymbol{\omega} \in \bar{\Sigma}^\infty$ which map to \mathbf{y} are sequences of the form $\boldsymbol{\omega} = \mathbf{y}\text{EOS}\boldsymbol{\omega}'$ for $\boldsymbol{\omega}' \in \bar{\Sigma}^\infty$ —this is exactly the cylinder $\bar{C}(\mathbf{y}\text{EOS})$! By the definition of the probability space $(\bar{\Sigma}, \sigma(\bar{C}), \mathbb{P})$, this is

$$\mathbb{P}(\bar{C}(\mathbf{y}\text{EOS})) = \sum_{\mathbf{y}' \in \{(\mathbf{y}\text{EOS})\}} p_{\text{LN}}(\mathbf{y}') = p_{\text{LN}}(\mathbf{y}\text{EOS}) \quad (2.49)$$

and as before $p_{\text{LN}}(\mathbf{y}\text{EOS}) = \prod_{t=1}^T p_{\text{LN}}(y_t \mid \mathbf{y}_{<t}) p_{\text{LN}}(\text{EOS} \mid \mathbf{y})$.

Altogether, this means that, given a finite string $\mathbf{y} \in \Sigma^*$, we intuitively have

$$\mathbb{P}_*(\mathbf{x} = \mathbf{y}) = p_{\text{LN}}(\text{EOS} \mid \mathbf{y}) p_{\text{LN}}(\mathbf{y}). \quad (2.50)$$

Additionally, as we will show in the next section, the probability of the set of infinite strings $\mathbb{P}_*(\mathbf{x} \in \Sigma^\infty)$ is the probability of generating an infinite string.

An important technical detail left out in this discussion so far is that both the singleton set $\{\mathbf{y}\}$ and Σ^∞ need to be measurable in $(\Sigma^* \cup \Sigma^\infty, \sigma(\mathcal{C}))$ for the above to make sense. This is addressed by Proposition 2.5.3 and Proposition 2.5.4.

Proposition 2.5.3. *In measure space $(\Sigma^* \cup \Sigma^\infty, \sigma(\mathcal{C}))$, $\{\mathbf{y}\}$ is measurable for all $\mathbf{y} \in \Sigma^*$.*

Proof. By definition in Eq. (2.43), for any $\mathbf{y} \in \Sigma^*$,

$$\mathcal{C}(\mathbf{y}) = \{\mathbf{y}\boldsymbol{\omega} : \boldsymbol{\omega} \in \Sigma^* \cup \Sigma^\infty\} \quad (2.51a)$$

$$= \{\mathbf{y}\boldsymbol{\omega} : \boldsymbol{\omega} \in \Sigma^*\} \cup \{\mathbf{y}\boldsymbol{\omega} : \boldsymbol{\omega} \in \Sigma^\infty\} \quad (2.51b)$$

where

$$\{\mathbf{y}\boldsymbol{\omega} : \boldsymbol{\omega} \in \Sigma^*\} = \{\mathbf{y}\} \cup \bigcup_{a \in \Sigma} \{\mathbf{y}a\boldsymbol{\omega} : \boldsymbol{\omega} \in \Sigma^*\} \quad (2.52a)$$

and

$$\{\mathbf{y}\boldsymbol{\omega} : \boldsymbol{\omega} \in \Sigma^\infty\} = \bigcup_{a \in \Sigma} \{\mathbf{y}a\boldsymbol{\omega} : \boldsymbol{\omega} \in \Sigma^\infty\}. \quad (2.53)$$

So,

$$\mathcal{C}(\mathbf{y}) = \{\mathbf{y}\} \cup \bigcup_{a \in \Sigma} \left(\{\mathbf{y}a\boldsymbol{\omega} : \boldsymbol{\omega} \in \Sigma^*\} \cup \{\mathbf{y}a\boldsymbol{\omega} : \boldsymbol{\omega} \in \Sigma^\infty\} \right) \quad (2.54a)$$

$$= \{\mathbf{y}\} \cup \bigcup_{a \in \Sigma} \mathcal{C}(\mathbf{y}a) \quad (2.54b)$$

which implies that $\{\mathbf{y}\} = \mathcal{C}(\mathbf{y}) \setminus \bigcup_{a \in \Sigma} \mathcal{C}(\mathbf{y}a)$ and hence measurable. ■

Proposition 2.5.4. *In the measure space $(\Sigma^* \cup \Sigma^\infty, \sigma(\mathcal{C}))$, Σ^∞ is measurable.*

Proof. First, the outcome space $\Sigma^* \cup \Sigma^\infty$ is measurable by definition of σ -algebra. Notice that

$$\Sigma^\infty = (\Sigma^* \cup \Sigma^\infty) \setminus \bigcup_{\mathbf{y} \in \Sigma^*} \{\mathbf{y}\}. \quad (2.55)$$

Since each $\{\mathbf{y}\}$ in the above is measurable by Proposition 2.5.3 and Σ^* is a countable set, Σ^∞ is then measurable. ■

Since both $\{\mathbf{y}\}$ and Σ^∞ are measurable in $(\Sigma^* \cup \Sigma^\infty, \sigma(\mathcal{C}))$ by Propositions 2.5.3 and 2.5.4, we have the following.

Proposition 2.5.5. *A sequence model $(\Sigma^* \cup \Sigma^\infty, \sigma(\mathcal{C}), \mathbb{P})$ is tight if and only if $\sum_{\mathbf{y} \in \Sigma^*} \mathbb{P}(\{\mathbf{y}\}) = 1$.*

Proof. By definition, a sequence model is tight if and only if $\mathbb{P}(\Sigma^\infty) = 0$. By Propositions 2.5.3 and 2.5.4, we can write

$$\mathbb{P}(\Sigma^* \cup \Sigma^\infty) = \mathbb{P}(\Sigma^\infty) + \mathbb{P}(\Sigma^*) \quad (\text{countable additivity}) \quad (2.56a)$$

$$= \mathbb{P}(\Sigma^\infty) + \sum_{\mathbf{y} \in \Sigma^*} \mathbb{P}(\{\mathbf{y}\}). \quad (\text{countable additivity}) \quad (2.56b)$$

Hence, a sequence model is tight if and only if $\sum_{\mathbf{y} \in \Sigma^*} \mathbb{P}(\{\mathbf{y}\}) = 1$. ■

Deriving eos

As an aside, the preceding section allows us to motivate the EOS token in LNM as a construct that emerges naturally. Specifically, for any $\mathbf{y} \in \Sigma^*$, rearranging Eq. (2.50):

$$p_{\text{LN}}(\text{EOS} \mid \mathbf{y}) = \frac{\mathbb{P}_*(\mathbf{x} = \mathbf{y})}{p_{\text{LN}}(\mathbf{y})} \quad (2.57a)$$

$$= \frac{\mathbb{P}_*(\mathbf{x} = \mathbf{y})}{\mathbb{P}_*(\mathbf{x} \in \mathcal{C}(\mathbf{y}))} \quad (2.57b)$$

$$= \mathbb{P}_*(\mathbf{x} = \mathbf{y} \mid \mathbf{x} \in \mathcal{C}(\mathbf{y})) \quad (2.57c)$$

where we have used $p_{\text{LN}}(\mathbf{y}) = \mathbb{P}(\overline{\mathcal{C}}(\mathbf{y})) = \mathbb{P}(\mathbf{x}^{-1}(\mathcal{C}(\mathbf{y}))) = \mathbb{P}_*(\mathbf{x} \in \mathcal{C}(\mathbf{y}))$. This means that the EOS probability in an LNM emerges as the conditional probability that, given that we must generate a string with a prefix $\mathbf{y} \in \Sigma^*$, the string is *exactly* \mathbf{y} , i.e., that generation ends there.

2.5.4 Characterizing Tightness

Now that we have derived a measure-theoretic formalization of the probability space induced by locally-normalized models, we can use it to provide an exact characterization of tightness in LNM. First, we consider the event

$$\mathcal{A}_k \stackrel{\text{def}}{=} \{\boldsymbol{\omega} \in \overline{\Sigma}^\infty : \omega_k = \text{EOS}\} \quad (2.58)$$

in the probability space $(\overline{\Sigma}^\infty, \sigma(\overline{\mathcal{C}}), \mathbb{P})$. Intuitively, \mathcal{A}_k is the event that an EOS symbol appears at position k in the string. Note that under this definition the \mathcal{A}_k are not disjoint. For example, the string $\boldsymbol{\omega} = ab\text{EOS}c\text{EOS}d\text{EOS}d\text{EOS}\dots$ lives in the intersection of \mathcal{A}_3 and \mathcal{A}_5 since EOS appears at both position 3 and position 5. Using Eq. (2.58), we can express the event consisting of all finite strings as

$$\bigcup_{k=1}^{\infty} \mathcal{A}_k. \quad (2.59)$$

It follows that we can express the event of an infinite string as

$$\left(\bigcup_{k=1}^{\infty} \mathcal{A}_k \right)^c = \bigcap_{k=1}^{\infty} \mathcal{A}_k^c. \quad (2.60)$$

Thus, using the random variable x , we can express the probability of generating an infinite string as

$$\mathbb{P}_*(x \in \Sigma^\infty) = \mathbb{P}(x^{-1}(\Sigma^\infty)) \quad (2.61a)$$

$$= \mathbb{P}\left(\bigcap_{k=1}^{\infty} \mathcal{A}_k^c\right). \quad (2.61b)$$

Hence, we can now restate and formalize the notion of tightness.

Definition 2.5.5. A sequence model is said to be **tight** if $\mathbb{P}_*(x \in \Sigma^\infty) = 0$, in which case it is also a language model. Otherwise, we say that it is **non-tight**.

Note that the definition of \mathcal{A}_k only uses a string's k -prefix, and hence is a cylinder set of rank k . Recalling that the cylinder sets are measurable and so are the sets countably generated by them, we see that both the event consisting of all finite strings and the event consisting of all infinite strings are measurable. Thus, $\mathbb{P}(\cup_{k=1}^{\infty} \mathcal{A}_k)$ and $\mathbb{P}(\cap_{k=1}^{\infty} \mathcal{A}_k^c)$ are well defined.

A Lower Bound Result

We have characterized tightness in terms of the probability of a specific event $\mathbb{P}(\cap_{k=1}^{\infty} \mathcal{A}_k^c)$, a quantity we now seek to determine.

Lemma 2.5.4. If $\sum_{n=2}^{\infty} \mathbb{P}(\mathcal{A}_n \mid \cap_{m=1}^{n-1} \mathcal{A}_m^c) = \infty$, then $\mathbb{P}(\cap_{m=1}^{\infty} \mathcal{A}_m^c) = 0$.

Proof. First, recall an elementary inequality that for $x > 0$,

$$x - 1 \geq \log x \quad \Leftrightarrow \quad 1 - x \leq \log \frac{1}{x}. \quad (2.62)$$

Note that $\mathbb{P}(\cap_{m=1}^n \mathcal{A}_m^c) > 0$ for any n , for otherwise the conditional probabilities would be undefined. Let $p_n \stackrel{\text{def}}{=} \mathbb{P}(\cap_{m=1}^n \mathcal{A}_m^c)$. Then we have that $p_n > 0$ for all n , and

$$\infty = \sum_{n=2}^{\infty} \mathbb{P}(\mathcal{A}_n \mid \cap_{m=1}^{n-1} \mathcal{A}_m^c) \quad (2.63a)$$

$$= \sum_{n=2}^{\infty} 1 - \mathbb{P}(\mathcal{A}_n^c \mid \cap_{m=1}^{n-1} \mathcal{A}_m^c) \quad (2.63b)$$

$$= \lim_{N \rightarrow \infty} \sum_{n=2}^N 1 - \mathbb{P}(\mathcal{A}_n^c \mid \cap_{m=1}^{n-1} \mathcal{A}_m^c) \quad (2.63c)$$

$$\leq \lim_{N \rightarrow \infty} \sum_{n=2}^N \log 1 / \mathbb{P}(\mathcal{A}_n^c \mid \cap_{m=1}^{n-1} \mathcal{A}_m^c) \quad (\text{by Eq. (2.62)}) \quad (2.63d)$$

$$= \lim_{N \rightarrow \infty} \sum_{n=2}^N \log \frac{\mathbb{P}(\cap_{m=1}^{n-1} \mathcal{A}_m^c)}{\mathbb{P}(\cap_{m=1}^n \mathcal{A}_m^c)} \quad (2.63e)$$

$$= \lim_{N \rightarrow \infty} \sum_{n=2}^N \log \frac{p_{n-1}}{p_n} \quad (2.63f)$$

$$= \lim_{N \rightarrow \infty} \sum_{n=2}^N (\log p_{n-1} - \log p_n) \quad (2.63g)$$

$$= \lim_{N \rightarrow \infty} (\log p_1 - \log p_N) \quad (2.63h)$$

$$= \log p_1 - \lim_{N \rightarrow \infty} \log p_N \quad (2.63i)$$

which implies that

$$\lim_{N \rightarrow \infty} \log p_N = -\infty \quad (2.64a)$$

$$\Leftrightarrow \lim_{N \rightarrow \infty} p_N = 0 \quad (2.64b)$$

$$\Leftrightarrow \lim_{N \rightarrow \infty} \mathbb{P}(\cap_{m=1}^N \mathcal{A}_m^c) = 0 \quad (2.64c)$$

$$\Leftrightarrow \mathbb{P}(\cap_{m=1}^\infty \mathcal{A}_m^c) = 0. \quad (\text{by continuity of measure}) \quad (2.64d)$$

■

Using Lemma 2.5.4, we can derive the following useful condition of tightness of a language model. Specifically, it applies when the probability of EOS is lower bounded by a function that depends only on the *length* and not the *content* of the prefix.

Proposition 2.5.6. *If $p_{LN}(\text{EOS} \mid \mathbf{y}) \geq f(t)$ for all $\mathbf{y} \in \Sigma^t$ and for all t and $\sum_{t=1}^\infty f(t) = \infty$, then $\mathbb{P}(\cap_{k=1}^\infty \mathcal{A}_k^c) = 0$. In other words, p_{LN} is tight.*

Proof. Suppose $p_{LN}(\text{EOS} \mid \mathbf{y}) \geq f(t)$ for all $\mathbf{y} \in \Sigma^t$. To apply Lemma 2.5.4, we observe that

$$\mathcal{A}_n \cap (\mathcal{A}_1^c \cap \cdots \cap \mathcal{A}_{n-1}^c) = \{\omega \in \bar{\Sigma}^\infty : \omega_n = \text{EOS}\} \cap \left(\bigcap_{i=1}^{n-1} \{\omega \in \bar{\Sigma}^\infty : \omega_i \neq \text{EOS}\} \right) \quad (2.65a)$$

$$= \{\omega \in \bar{\Sigma}^\infty : \omega = \text{EOS}, \forall i < n, \omega \neq \text{EOS}\} \quad (2.65b)$$

$$= \{\omega \in \bar{\Sigma}^\infty : \omega \text{'s first EOS is at position } n\} \quad (2.65c)$$

and similarly

$$\mathcal{A}_1^c \cap \cdots \cap \mathcal{A}_{n-1}^c = \{\omega \in \bar{\Sigma}^\infty : \text{There is no EOS in } \omega \text{'s first } n-1 \text{ positions}\} \quad (2.66)$$

Setting $\mathcal{G} \stackrel{\text{def}}{=} \{\omega \text{EOS} : \omega \in \Sigma^{n-1}\} \subset \bar{\Sigma}^n$, we get

$$\mathbb{P}(\mathcal{A}_n \mid \mathcal{A}_1^c \cap \cdots \cap \mathcal{A}_{n-1}^c) = \frac{\mathbb{P}(\mathcal{A}_n \cap (\mathcal{A}_1^c \cap \cdots \cap \mathcal{A}_{n-1}^c))}{\mathbb{P}(\mathcal{A}_1^c \cap \cdots \cap \mathcal{A}_{n-1}^c)} \quad (2.67a)$$

$$= \frac{\mathbb{P}(\bar{C}(\mathcal{G}))}{\mathbb{P}(\bar{C}(\Sigma^{n-1}))} \quad (\text{definition of } \mathcal{G}) \quad (2.67b)$$

$$= \frac{\sum_{\omega \in \Sigma^{n-1}} p(\text{EOS} \mid \omega) p(\omega)}{\sum_{\omega \in \Sigma^{n-1}} p(\omega)} \quad (\text{by Eq. (2.30)}) \quad (2.67c)$$

$$\geq \frac{\sum_{\omega \in \Sigma^{n-1}} f(n-1) p(\omega)}{\sum_{\omega \in \Sigma^{n-1}} p(\omega)} \quad (\text{definition of } f(t)) \quad (2.67d)$$

$$= f(n-1) \frac{\sum_{\omega \in \Sigma^{n-1}} p(\omega)}{\sum_{\omega \in \Sigma^{n-1}} p(\omega)} \quad (2.67e)$$

$$= f(n-1). \quad (2.67f)$$

Since $\sum_{t=0}^\infty f(t) = \infty$, Lemma 2.5.4 shows that the event of a string never terminating, i.e., $\cap_{k=1}^\infty \mathcal{A}_k^c$ has probability measure $\mathbb{P}(\cap_{k=1}^\infty \mathcal{A}_k^c) = 0$. In other words, if the EOS probability of a language model is lower bounded by a divergent sequence at every step, then the event that this language model terminates has probability 1. ■

The Borel–Cantelli Lemmata

It turns out that Proposition 2.5.6 admits a converse statement in which we can prove a similar property of p_{LN} by assuming that the model is tight. To show this result, we will use a fundamental inequality from probability theory—the Borel–Cantelli lemmata. The Borel–Cantelli lemmata are useful for our purposes because they relate the probability measure of sets of the form $\bigcap_{n=0}^\infty \mathcal{A}_n$ or $\bigcup_{n=0}^\infty \mathcal{A}_n$ to a series $\sum_{n=0}^\infty p_n$. We

will only state the lemmata here without supplying their proofs;²⁴ however, we point out that Lemma 2.5.4 can be viewed as a parallel statement to the Borel–Cantelli lemmata and one can prove the lemmata using a very similar proof (cf. proof of Theorem 2.3.7 in [Durrett, 2019](#)).

Concretely, given a sequence of events $\{\mathcal{A}_n\}_{n=1}^\infty$ in some probability space, the Borel–Cantelli lemmata are statements about the event

$$\{\mathcal{A}_n \text{ i.o.}\} \stackrel{\text{def}}{=} \bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} \mathcal{A}_n \quad (2.68)$$

where i.o. stands for “infinitely often.” Intuitively, $\{\mathcal{A}_n \text{ i.o.}\}$ is the set of outcomes that appear in infinitely many sets in the collection $\{\mathcal{A}_n\}_{n=1}^\infty$ —they are the events that always remain in the union of an infinite family of sets no matter how many of the leading ones we remove (hence the name). We will not use Borel–Cantelli directly, but they offer a probabilistic proof of a key result (Corollary 2.5.1) which will in turn lead to the desired statement about tightness. We formally state the first and second Borel–Cantelli lemmata below.

Lemma 2.5.5 (Borel–Cantelli I). *If $\sum_{n=1}^\infty \mathbb{P}(\mathcal{A}_n) < \infty$, then $\mathbb{P}(\mathcal{A}_n \text{ i.o.}) = 0$.*

Lemma 2.5.6 (Borel–Cantelli II). *Assume $\{\mathcal{A}_n\}$ is a sequence of independent events, then $\sum_{n=1}^\infty \mathbb{P}(\mathcal{A}_n) = \infty \Rightarrow \mathbb{P}(\mathcal{A}_n \text{ i.o.}) = 1$.*

Using the Borel–Cantelli lemmata, we can prove the following useful fact.

Corollary 2.5.1. *Given a sequence $\{p_n\}$ where $p_n \in [0, 1)$. Then,*

$$\prod_{n=1}^{\infty} (1 - p_n) = 0 \iff \sum_{n=1}^{\infty} p_n = \infty. \quad (2.69)$$

To show Corollary 2.5.1, we first show the following simple consequence of Borel–Cantelli.

Corollary 2.5.2. *If $\mathbb{P}(\mathcal{A}_n \text{ i.o.}) = 1$, then $\sum_{n=1}^\infty \mathbb{P}(\mathcal{A}_n) = \infty$.*

Proof. Suppose to the contrary that $\sum_{n=1}^\infty \mathbb{P}(\mathcal{A}_n) < \infty$, then, by Borel–Cantelli I (Lemma 2.5.5), $\mathbb{P}(\mathcal{A}_n \text{ i.o.}) = 0$, which contradicts the assumption. Hence, $\sum_{n=1}^\infty \mathbb{P}(\mathcal{A}_n) = \infty$. ■

Proof. We can use a product measure to construct a sequence of independent events $\{\mathcal{A}_n\}_{n=1}^\infty$ such that $\mathbb{P}(\mathcal{A}_n) = p_n$. (The product measure ensures independence.) Then, by definition in Eq. (2.68),

$$\{\mathcal{A}_n \text{ i.o.}\}^c = \bigcup_{m=1}^{\infty} \bigcap_{n \geq m} \mathcal{A}_n^c \quad (2.70)$$

So,

$$1 - \mathbb{P}(\mathcal{A}_n \text{ i.o.}) = \mathbb{P}\left(\bigcup_m \bigcap_{n \geq m} \mathcal{A}_n^c\right) \quad (2.71a)$$

$$= \lim_{m \rightarrow \infty} \mathbb{P}\left(\bigcap_{n \geq m} \mathcal{A}_n^c\right) \quad (2.71b)$$

$$= \lim_{m \rightarrow \infty} \prod_{n \geq m} \mathbb{P}(\mathcal{A}_n^c) \quad (\mathcal{A}_n \text{ are independent by construction}) \quad (2.71c)$$

$$= \lim_{m \rightarrow \infty} \prod_{n \geq m} (1 - p_n) \quad (2.71d)$$

²⁴See §2.3 in [Durrett \(2019\)](#) or §4 in [Billingsley \(1995\)](#) instead.

(\Rightarrow): Assume $\prod_{n=1}^{\infty} (1 - p_n) = 0$. Then, for any m ,

$$0 = \prod_{n \geq 1} (1 - p_n) = \underbrace{\left(\prod_{1 \leq n < m} (1 - p_n) \right)}_{>0} \left(\prod_{n \geq m} (1 - p_n) \right) \quad (2.72)$$

So it must be the case that, for any m , $\prod_{n \geq m} (1 - p_n) = 0$. Therefore,

$$1 - \mathbb{P}(\mathcal{A}_n \text{ i.o.}) = \lim_{m \rightarrow \infty} \prod_{n \geq m} (1 - p_n) = 0 \quad (2.73)$$

which implies $\mathbb{P}(\mathcal{A}_n \text{ i.o.}) = 1$. Corollary 2.5.2 implies that $\sum_{n=1}^{\infty} p_n = \infty$.

(\Leftarrow): Assume $\sum_{n=1}^{\infty} p_n = \infty$. Then by Borel–Cantelli II (Lemma 2.5.6), $\mathbb{P}(\mathcal{A}_n \text{ i.o.}) = 1$ which implies

$$0 = 1 - \mathbb{P}(\mathcal{A}_n \text{ i.o.}) = \lim_{m \rightarrow \infty} \prod_{n \geq m} (1 - p_n) \quad (2.74)$$

Observe that $\left\{ \prod_{n \geq m} (1 - p_n) \right\}_m$ is a non-increasing sequence in m ; to see this, note that as m grows larger we multiply strictly fewer values $(1 - p_n) \in (0, 1]$. However, since we know the sequence is non-negative and tends to 0, it follows that for *any* m , we have

$$\prod_{n \geq m} (1 - p_n) = 0. \quad (2.75)$$

It follows that, for any m , we have

$$\prod_{n=1}^{\infty} (1 - p_n) = \prod_{n < m} (1 - p_n) \underbrace{\prod_{n \geq m} (1 - p_n)}_{=0} = \prod_{n < m} (1 - p_n) \cdot 0 = 0. \quad (2.76)$$

■

We now turn to proving a more general version of Proposition 2.5.6, which would imply its converse. First, we define the following quantity

$$\tilde{p}_{\text{EOS}}(t) \stackrel{\text{def}}{=} \mathbb{P}(\mathcal{A}_t \mid \mathcal{A}_1^c \cap \cdots \cap \mathcal{A}_{t-1}^c) \quad (2.77)$$

which can be viewed as the EOS probability at step t , given that EOS was not generated at any earlier step. One can also show that, when $\tilde{p}_{\text{EOS}}(t)$ is defined, it has the same value as

$$\tilde{p}_{\text{EOS}}(t) = \frac{\sum_{\omega \in \Sigma^{t-1}} p_{\text{LN}}(\omega) p_{\text{LN}}(\text{EOS} \mid \omega)}{\sum_{\omega \in \Sigma^{t-1}} p_{\text{LN}}(\omega)}, \quad (2.78)$$

which one can see as the weighted average probability of terminating at a string of length t .

We can now completely characterize the tightness of an LNM with the following theorem.

Theorem 2.5.3. *An LNM is tight if and only if $\tilde{p}_{\text{EOS}}(t) = 1$ for some t or $\sum_{t=1}^{\infty} \tilde{p}_{\text{EOS}}(t) = \infty$.*

Proof. Recall the definition of \tilde{p}_{EOS} , as previously defined in Eq. (2.77), is

$$\tilde{p}_{\text{EOS}}(t) \stackrel{\text{def}}{=} \mathbb{P}(\mathcal{A}_t \mid \mathcal{A}_1^c \cap \cdots \cap \mathcal{A}_{t-1}^c). \quad (2.79)$$

Case 1. Suppose that $\tilde{p}_{\text{EOS}}(t) < 1$ for all t . Consider the termination probability again:

$$\mathbb{P}\left(\bigcap_{t=1}^{\infty} \mathcal{A}_t^c\right) = \lim_{T \rightarrow \infty} \mathbb{P}\left(\bigcap_{t=1}^T \mathcal{A}_t^c\right) \quad (2.80a)$$

$$= \lim_{T \rightarrow \infty} \prod_{t=1}^T \mathbb{P}(\mathcal{A}_t^c \mid \mathcal{A}_1^c \cap \dots \cap \mathcal{A}_{t-1}^c) \quad (2.80b)$$

$$= \lim_{T \rightarrow \infty} \prod_{t=1}^T (1 - \tilde{p}_{\text{EOS}}(t)) \quad (2.80c)$$

$$= \prod_{t=1}^{\infty} (1 - \tilde{p}_{\text{EOS}}(t)). \quad (2.80d)$$

In the above, we have assumed that $\mathbb{P}(\mathcal{A}_1^c \cap \dots \cap \mathcal{A}_t^c) > 0$ for all t , which is true by assumption that $\tilde{p}_{\text{EOS}}(t) < 1$. Hence, by Corollary 2.5.1, Eq. (2.80d) is 0 if and only if $\sum_t \tilde{p}_{\text{EOS}}(t) = \infty$.

Case 2. If $\tilde{p}_{\text{EOS}}(t) = 1$ is true for some $t = t_0$, then $\mathbb{P}(\mathcal{A}_1^c \cap \dots \cap \mathcal{A}_{t_0}^c) = 0$ and hence $\mathbb{P}(\bigcap_{t=1}^{\infty} \mathcal{A}_t^c) = 0$ and such a language model is guaranteed to terminate at t_0 . ■

The first condition intuitively says that there exists a step t at which the LNM will stop with probability 1. If the first case of the condition does not hold, the second case can be checked since its summands will be well-defined (the conditional probabilities in Eq. (2.78) will not divide by 0). We remark that Theorem 2.5.3 is a generalization of Proposition 2.5.6 since if $\tilde{p}_{\text{EOS}}(t)$ is lower-bounded by $f(t)$ whose series diverges, its own series would also diverge. However, since $\tilde{p}_{\text{EOS}}(t)$ involves the computation of a partition function in its denominator, it is most likely intractable to calculate (Lin et al., 2021a; Lin and McCarthy, 2022). Hence, Proposition 2.5.6 will be the main tool for determining tightness when we explore concrete language modeling frameworks later.

We have now very thoroughly defined the notion of language model tightness and provided sufficient and necessary conditions for an LNM or a sequence model to be tight. In the next sections, we start our exploration of concrete computational models of language, from the very simple and historically important finite-state language models, their neural variants, to the modern Transformer architectures. For each of them, we will also individually discuss their tightness results and conditions.

Chapter 3

Modeling Foundations

The previous chapter introduced the fundamental measure-theoretic characteristics of language modeling. We will revisit those over and over as they will serve as the foundations on which subsequent concepts are built.

In this chapter, we turn our attention to *modeling* foundations, that is, the decisions we face when we want to *build* a distribution over strings and *learn* the appropriate parameters for that distribution. We first discuss *how* to parameterize a distribution over strings (§3.1), what it means to *learn* good parameters, and how this can be done with modern optimization techniques and objectives (§3.2).

Continuing our framing of the notes in terms of questions, we will try to address the following:

Question 3.0.1. *How can a sequence model be parameterized?*

We introduce a more formal definition of a “parametrized model” later. For now, you can simply think of it as a function $p_{\theta} : \bar{\Sigma}^* \rightarrow \mathbb{R}$ described by some *free parameters* $\theta \in \Theta$ from a parameter space Θ . This means that the values that p_{θ} maps its inputs to might depend on the choice of the parameters θ —the presence of parameters in a model, therefore, allows us to *fit* them, which in our context specifically, means choosing them to maximize some objective with respect to data. This raises the following question:

Question 3.0.2. *Given a parameterized model and a dataset, how can model parameters be chosen to reflect the dataset as well as possible?*

We begin with Question 3.0.1.

3.1 Representation-based Language Models

Most modern language models are defined as locally normalized models. However, in order to define locally normalized language model, we first define a sequence model $p_{\text{SM}}(\bar{y} | \bar{y})$. Then, we prove that the specific parameterization used in $p_{\text{SM}}(\bar{y} | \bar{y})$ encodes a tight locally normalized language model. However, as we demonstrated in Example 2.5.1, not all sequence models encode tight locally normalized language models in the sense of definition Definition 2.5.1. So far, however, we have only talked about this process abstractly. For example, we have proven that every language model can be locally normalized and we have also given necessary and sufficient conditions for when a sequence model encodes a tight locally normalized language model. In this section, we start making the abstraction more concrete by considering a very general framework for parameterizing a locally normalized language model through sequence models $p_{\text{SM}}(\bar{y} | \bar{y})$. We will call this the **representation-based language modeling** framework.

In the representation-based language modeling framework, each conditional distribution in a sequence model $p_{\text{SM}}(\bar{y} | \bar{y})$ directly models the probability of the next symbol $\bar{y} \in \Sigma$ given the context \bar{y} —in other words, it tells us how likely \bar{y} is appears in the context of \bar{y} .¹ For example, given the string $\bar{y} = \text{“Papa eats caviar with a”}$, we would like $p_{\text{LN}}(\bar{y} | \bar{y})$ to capture that “spoon” is more likely than “fork”. At the same time, since eating caviar with a fork is technically possible, we would also like $p_{\text{LN}}(\bar{y} | \bar{y})$ to capture that “fork” is likelier than, for example, “pencil”.

However, it is not *a-priori* clear how we should model $p_{\text{SM}}(\bar{y} | \bar{y})$ concretely. We want to define a function that can map contexts \bar{y} to a distribution over possible continuations \bar{y} with the caveat that this distribution can be easily adjusted, i.e., we can optimize its parameters with some objective in mind (cf. §3.2). We will do this by adopting a very general idea of defining $p_{\text{SM}}(\bar{y} | \bar{y})$ in terms of similarity between representations that represent the symbol \bar{y} and the context \bar{y} . The more compatible the symbol \bar{y} is with the context \bar{y} , the more probable it should be. Intuitively, going from the example above, this means that “spoon” should be more similar to “Papa eats caviar with a” than “fork” should be, and that should still be more similar than “pencil”. On the other hand, notice that this also means that “spoon” and “fork” should be closer together than any of them to “pencil”.

One possibility for doing this is by *embedding* individual symbols \bar{y} and all possible contexts \bar{y} as vectors in a Hilbert space, i.e., a complete vector space endowed with an inner product. Once we embed the symbols and contexts in such a space, we can talk about how similar they are. We will first describe how this can be done abstractly §3.1.1 and then discuss how exactly vector representations can be used when defining *discrete* probability distributions over the symbols in §3.1.3 by taking into account the notion of similarities between vectors. We discuss methods for *learning* representations later in this chapter (§3.2) and in Chapter 5.

3.1.1 Vector Space Representations

It is not immediately obvious how to measure the similarity or compatibility between two symbols, two contexts or a symbol and a context. However, such a notion is required as part of our intuitive desiderata for $p_{\text{SM}}(\bar{y} | \bar{y})$. We begin by stating an important guiding principle, which we describe in detail next and use heavily throughout the rest of the notes.

Principle 3.1.1 (Representation Learning). *The **good representation principle** states that the success of a machine learning model depends—in great part—on the representation that is chosen (or learned) for the objects that are being modeled. In the case of language modeling, the two most salient choice points are the representations chosen for the symbols, elements of Σ , and the representations chosen for the contexts, elements of Σ^* .*

Learning vector representations from data where individual entities are represented in some **representation space** (i.e., a Hilbert space) has a rich history in NLP and machine learning in general (Bengio et al., 2013).

To discuss the representations of symbols and strings more formally, we first introduce the notion of a **Hilbert space**, which leads us to a useful geometric manner to discuss the similarity and compatibility of

¹Unless explicitly stated otherwise, we use the phrase “in the context of” to imply given *prior* context—i.e., when discussing probability distributions, this refers to the distribution $p_{\text{SM}}(\bar{y}_t | \bar{y}_{<t})$ with $\bar{y} = \bar{y}_t$. We will also see examples of models which specify the conditional probabilities in terms of symbols that do not necessarily appear before the current one.

symbols and contexts. We first start with some more basic definitions. A vector space over a field \mathbb{F} is a set \mathbb{V} together with two binary operations that satisfy certain axioms. The elements of \mathbb{F} are often referred to as **scalars** and the elements of \mathbb{V} as **vectors**. The two operations in the definition of a vector space are the *addition of vectors* and *scalar multiplication of vectors*.

Definition 3.1.1 (Vector space). *A **vector space** over a field \mathbb{F} is a set \mathbb{V} together with two binary operations that satisfy the following axioms:*

1. **Associativity** of vector addition: for all $\mathbf{v}, \mathbf{u}, \mathbf{q} \in \mathbb{V}$

$$(\mathbf{v} + \mathbf{u}) + \mathbf{q} = \mathbf{v} + (\mathbf{u} + \mathbf{q}) \quad (3.1)$$

2. **Commutativity** of vector addition: for all $\mathbf{v}, \mathbf{u} \in \mathbb{V}$

$$\mathbf{v} + \mathbf{u} = \mathbf{u} + \mathbf{v} \quad (3.2)$$

3. **Identity** element of vector addition: there exists $\mathbf{0} \in \mathbb{V}$ such that for all $\mathbf{v} \in \mathbb{V}$

$$\mathbf{v} + \mathbf{0} = \mathbf{v} \quad (3.3)$$

4. **Inverse** elements of vector addition: for every $\mathbf{v} \in \mathbb{V}$ there exists a $-\mathbf{v} \in \mathbb{V}$ such that

$$\mathbf{v} + (-\mathbf{v}) = \mathbf{0} \quad (3.4)$$

5. **Compatibility** of scalar multiplication with field multiplication: for all $\mathbf{v} \in \mathbb{V}$ and $x, y \in \mathbb{F}$

$$x(y\mathbf{v}) = (xy)\mathbf{v} \quad (3.5)$$

6. **Identity** element of scalar multiplication: for all $\mathbf{v} \in \mathbb{V}$

$$1\mathbf{v} = \mathbf{v} \quad (3.6)$$

where 1 is the multiplicative identity in \mathbb{F} .

7. **Distributivity** of scalar multiplication with respect to vector addition: for all $x \in \mathbb{F}$ and all $\mathbf{u}, \mathbf{v} \in \mathbb{V}$

$$x(\mathbf{v} + \mathbf{u}) = x\mathbf{v} + x\mathbf{u} \quad (3.7)$$

8. **Distributivity** of scalar multiplication with respect to field addition: for all $x, y \in \mathbb{F}$ and all $\mathbf{v} \in \mathbb{V}$

$$(x + y)\mathbf{v} = x\mathbf{v} + y\mathbf{v} \quad (3.8)$$

In almost all practical cases, \mathbb{F} will be \mathbb{R} and \mathbb{V} will be \mathbb{R}^D for some $D \in \mathbb{N}$.

An important characteristic of a vector space is its **dimensionality**, which, informally, corresponds to the number of independent directions—**basis vectors**—in the space. Any $\mathbf{v} \in \mathbb{V}$ can be expressed as a *linear combination* of the D basis vectors. The coefficients of this linear combination can then be combined into a D -dimensional **coordinate vector** in \mathbb{F}^D . Vector spaces, therefore, allow us to talk about their elements in terms of their expressions with respect to the basis vectors. Inner product spaces additionally define an **inner product**, mapping pairs of elements of the vector space to scalars. More formally, it is a vector space together with a map $\langle \cdot, \cdot \rangle$ (the inner product) defined as follows.

Definition 3.1.2 (Inner product space). *An **inner product space** is a vector space \mathbb{V} over a field \mathbb{F} coupled with a map*

$$\langle \cdot, \cdot \rangle : \mathbb{V} \times \mathbb{V} \rightarrow \mathbb{F} \quad (3.9)$$

such that the following axioms hold

1. **Conjugate symmetry:** for all $\mathbf{v}, \mathbf{u} \in \mathbb{V}$

$$\langle \mathbf{v}, \mathbf{u} \rangle = \overline{\langle \mathbf{u}, \mathbf{v} \rangle} \quad (3.10)$$

where \bar{x} denotes the **conjugate** of the element $x \in \mathbb{F}$.

2. **Linearity** in the first argument: for all $\mathbf{v}, \mathbf{u}, \mathbf{z} \in \mathbb{V}$ and $x, y \in \mathbb{F}$

$$\langle x\mathbf{v} + y\mathbf{u}, \mathbf{z} \rangle = x\langle \mathbf{v}, \mathbf{z} \rangle + y\langle \mathbf{u}, \mathbf{z} \rangle \quad (3.11)$$

3. **Positive-definiteness:** for all $\mathbf{v} \neq 0$

$$\langle \mathbf{v}, \mathbf{v} \rangle > 0 \quad (3.12)$$

Inner products are often defined such that they capture some notion of similarity of the vectors in \mathbb{V} . We will use this when formally defining $p_{\text{SM}}(\bar{y} \mid \bar{y})$ in §3.1.2.

Every inner product on a real or complex vector space induces a vector norm defined as follows.

Definition 3.1.3 (Norm). *Given a vector space \mathbb{V} over \mathbb{R} or \mathbb{C} and an inner product $\langle \cdot, \cdot \rangle$ over it, the **norm** induced by the inner product is defined as the function $\|\cdot\| : \mathbb{V} \rightarrow \mathbb{R}_{\geq 0}$ where*

$$\|\mathbf{v}\| \stackrel{\text{def}}{=} \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle}. \quad (3.13)$$

A Hilbert space is then an inner product space in which all sequences of elements satisfy a useful property with respect to the norm defined by the inner product: every convergent series with respect to the norm converges to a vector in \mathbb{V} .

Definition 3.1.4 (Hilbert space). *A **Hilbert space** is an inner product space that is **complete** with respect to the norm defined by the inner product. An inner product space is complete with respect to the norm if every Cauchy sequence (an absolutely convergent sequence, i.e., a sequence whose elements become arbitrarily close to each other) converges to an element in \mathbb{V} . More precisely, an inner product space is complete if, for every series*

$$\sum_{n=1}^{\infty} \mathbf{v}_n \quad (3.14)$$

such that

$$\sum_{n=1}^{\infty} \|\mathbf{v}_n\| < \infty, \quad (3.15)$$

it holds that

$$\sum_{n=1}^{\infty} \mathbf{v}_n \in \mathbb{V}. \quad (3.16)$$

Note that even if an inner product space \mathbb{V} is not necessarily a Hilbert space, \mathbb{V} can always be *completed* to a Hilbert space.

Theorem 3.1.1 (Completion theorem for inner product spaces). *Any inner product space can be completed into a Hilbert space.*

We omit the proof for this theorem. More precisely, the inner product space can be completed into a Hilbert space by completing it with respect to the norm induced by the inner product on the space. For this reason, inner product spaces are also called pre-Hilbert spaces.

To motivate our slightly more elaborate treatment of representation spaces, we consider an example of a model which falls under our definition of a representation-based language model but would be ill-defined if it worked under any space with fewer axioms than a Hilbert space.

Space	Utility
Vector space	A space in which representations of symbols and string live. It also allows the expression of the vector representations in terms of the basis vectors.
Inner product space	Defines an inner product, which defines a norm and can measure similarity.
Hilbert space	There are no “holes” in the representation space with respect to the defined norm, since all convergent sequences converge into \mathbb{V} .

Table 3.1: The utility of different spaces introduced in this section.

Example 3.1.1 (A series of representations). *Recurrent neural networks are a type of neural network that sequentially process their input and compute the output (context representation) at time step t based on the output at time step $t - 1$: $\mathbf{h}_t = \mathbf{f}(\mathbf{h}_{t-1}, y_t)$. A formal definition of a recurrent neural network, which we subsequently in §5.1.2, is not required at the moment. However, note that a recurrent neural network with one-dimensional representations h could, for example, take the specific form*

$$h_t = \frac{1}{2}h_{t-1} + \frac{1}{h_{t-1}} \quad (3.17)$$

with $h_0 = 2$.

Suppose we chose the inner product space \mathbb{Q} over the field \mathbb{Q} for our representation space. All elements of the sequence h_t are indeed rational numbers. However, the limit of the sequence, which can be shown to be $\sqrt{2}$, is not in the inner product space! This shows that \mathbb{Q} is not a Hilbert space and that we must, in full generality, work with Hilbert spaces whenever we are dealing with possibly infinite sequences of data. The reason this is especially relevant for language modeling is the need to consider arbitrarily long strings (contexts), whose representations we would like to construct in a way similar to Eq. (3.17). Such representations can, therefore, approach a limiting representation outside the space whenever the representation space does not satisfy the axioms of a Hilbert space.

A summary of the utilities of the three algebraic spaces introduced in this subsection is summarized in Tab. 3.1.

Representation Functions

We can now introduce the notion of a general representation function.

Definition 3.1.5 (Representation function). *Let \mathcal{S} be a set and \mathbb{V} a Hilbert space over some field \mathbb{F} . A **representation function** \mathbf{f} for the elements of \mathcal{S} is a function of the form*

$$\mathbf{f}: \mathcal{S} \mapsto \mathbb{V}. \quad (3.18)$$

The dimensionality of the Hilbert space of the representations, D , is determined by the modeler. In NLP, D usually ranges between 10 to 10000.

Importantly, in the case that \mathcal{S} is finite, we can represent a representation function as a *matrix* $\mathbf{E} \in \mathbb{R}^{|\mathcal{S}| \times D}$ (assuming $\mathbb{V} = \mathbb{R}^D$ where the n^{th} row corresponds to the representation of the n^{th} element of \mathcal{S} . This method for representing \mathbf{f} is both more concise and will be useful for integrating the symbol representation function into a model, where matrix multiplications are often the most efficient way to implement such functions on modern hardware.

This is the case for the representations of the individual symbols \bar{y} from $\bar{\Sigma}$, where the representation function, which we will denote as $\mathbf{e}(\cdot)$, is implemented as a lookup into the embedding matrix $\mathbf{E} \in \mathbb{R}^{|\bar{\Sigma}| \times D}$, i.e., $\mathbf{e}(\bar{y}) = \mathbf{E}_{\bar{y}}$.² In this case, we will also refer to $\mathbf{e}(\cdot)$ as the embedding function.

²Here, we use the notation $\mathbf{E}_{\bar{y}}$ to refer to the lookup of the row in \mathbf{E} corresponding to \bar{y} .

Definition 3.1.6 (Symbol embedding function). Let Σ be an alphabet. An **embedding function** $\mathbf{e}(\cdot) : \bar{\Sigma} \rightarrow \mathbb{R}^D$ is a representation function of individual symbols $\bar{y} \in \bar{\Sigma}$.

The representations $\mathbf{e}(\bar{y})$ are commonly referred to as **embeddings**, but, for consistency, we will almost exclusively use the term representations in this text. Let us first consider possibly the simplest way to represent discrete symbols with real-valued vectors: **one-hot encodings**.

Example 3.1.2 (One-hot encodings). Let $n : \bar{\Sigma} \rightarrow \{1, \dots, |\bar{\Sigma}|\}$ be a bijection (i.e., an ordering of the alphabet, assigning an index to each symbol in Σ). A one-hot encoding $\llbracket \cdot \rrbracket$ is a representation function which assigns the symbol $\bar{y} \in \bar{\Sigma}$ the $n(\bar{y})^{\text{th}}$ basis vector:

$$\llbracket y \rrbracket \stackrel{\text{def}}{=} \mathbf{d}_{n(y)}, \quad (3.19)$$

where here \mathbf{d}_n is the n^{th} canonical basis vector, i.e., a vector of zeros with a 1 at position n .

While one-hot encodings are an easy way to create vector representations of symbols, they have a number of drawbacks. First, these representations are relatively large—we have $D = |\bar{\Sigma}|$ —and *sparse*, since only one of the dimensions is non-zero. Second, such representations are not ideal for capturing the variation in the *similarity* between different words. For example, the cosine similarity—a metric we will motivate in the next section for measuring the similarity between symbol representations—between symbols’ one-hot encodings is zero for all non-identical symbols. Ideally, we would like symbol representations to encode semantic information, in which case, a metric such as cosine similarity could be used to quantify semantic similarity. This motivates the use of more complex representation functions, which we subsequently discuss.

While most systems use this standard way of defining individual symbol representations using the embedding matrix, the way that the *context* is encoded (and what even is considered as context) is really the major difference between the different architectures which we will consider later. Naturally, since the set of all contexts is infinite, we cannot simply represent the representation function with a matrix. Rather, we define the representation of a context \bar{y} through an encoding function.

Definition 3.1.7 (Context encoding function). Let Σ be an alphabet. A **context encoding function** $\text{enc}(\cdot) : \bar{\Sigma}^* \rightarrow \mathbb{R}^D$ is a representation function of strings $\bar{y} \in \bar{\Sigma}^*$.³

We will refer to $\text{enc}(\bar{y})$ as the **encoding** of $\bar{y} \in \bar{\Sigma}^*$. In the general framework, we can simply consider the encoding function enc to be a black box—however, a major part of Chapter 5 will concern defining specific functions enc and analyzing their properties.

With this, we now know how we can represent the discrete symbols and histories as real-valued vectors. We next consider how to use such representations for defining probability distributions over the next symbol.

3.1.2 Compatibility of Symbol and Context

Inner products naturally give rise to the geometric notion of angle, by giving us the means to measure the similarity between two representations. Concretely, the smaller the angle between the two representations is, the more similar the two representations are. In a Hilbert space, we *define* the cosine of the angle θ between the two representations

$$\cos(\theta) \stackrel{\text{def}}{=} \frac{\langle \mathbf{u}, \mathbf{v} \rangle}{\|\mathbf{u}\| \|\mathbf{v}\|}. \quad (3.20)$$

The Cauchy–Schwartz inequality immediately gives us that $\cos(\theta) \in [-1, 1]$ since $-\|\mathbf{u}\| \|\mathbf{v}\| \leq \langle \mathbf{u}, \mathbf{v} \rangle \leq \|\mathbf{u}\| \|\mathbf{v}\|$. Traditionally, however, we take the *unnormalized* cosine similarity as our measure of similarity, which simply corresponds to the inner product of the Hilbert space.

Given a context representation $\text{enc}(\bar{y})$, we can compute its inner products with all symbol representations $\mathbf{e}(\bar{y})$:

$$\langle \mathbf{e}(\bar{y}), \text{enc}(\bar{y}) \rangle. \quad (3.21)$$

³Note that, to be completely consistent, the encoding function should be defined over the set $(\bar{\Sigma} \cup \{\text{BOS}\})^*$ to allow for the case when $y_0 = \text{BOS}$. However, unlike EOS, we do not necessarily require BOS in any formal setting, which is why we leave it out. We apologize for this inconsistency.

which can be achieved simply with a matrix-vector product:

$$\mathbf{E} \text{enc}(\bar{\mathbf{y}}). \quad (3.22)$$

$\mathbf{E} \text{enc}(\bar{\mathbf{y}}) \in \mathbb{R}^{|\Sigma|}$, therefore, has the nice property that each of the individual entries corresponds to the similarities of a particular symbol to the context $\bar{\mathbf{y}}$. For reasons that will become clear soon, the entries of the vector $\mathbf{E} \text{enc}(\bar{\mathbf{y}})$ are often called **scores** or **logits**. This brings us almost to the final formulation of the probability distribution $p_{\text{SM}}(\bar{\mathbf{y}} \mid \bar{\mathbf{y}})$.

If $\mathbf{E} \text{enc}(\bar{\mathbf{y}})$ encodes similarity or compatibility, then a natural way to model the probability distribution $p_{\text{SM}}(\bar{\mathbf{y}} \mid \bar{\mathbf{y}})$ would be as proportional to the inner product between $\mathbf{e}(\bar{\mathbf{y}})$ and $\text{enc}(\bar{\mathbf{y}})$. However, the inner product $\langle \mathbf{e}(\bar{\mathbf{y}}), \text{enc}(\bar{\mathbf{y}}) \rangle$ may be negative; further, the sum over the similarity between a context and all tokens is not necessarily 1. To resolve this, we have to introduce the last piece of the puzzle: transforming $\mathbf{E} \text{enc}(\bar{\mathbf{y}})$ into a valid discrete probability distribution by using a projection function.

3.1.3 Projecting onto the Simplex

In the previous subsections we discussed how to encode symbols and contexts in a Hilbert space and how an inner product gives us a natural notation of similarity between a potentially infinite number of items. We can now finally discuss how to create the conditional distribution $p_{\text{SM}}(\bar{\mathbf{y}} \mid \bar{\mathbf{y}})$, i.e., how we can *map* the real-valued $\mathbf{E} \text{enc}(\bar{\mathbf{y}})$ that encodes symbol–context similarities to a valid probability distribution—a vector on the **probability simplex**.

Projection Functions: Mapping Vectors onto the Probability Simplex

$p_{\text{SM}}(\bar{\mathbf{y}} \mid \bar{\mathbf{y}})$ is a **categorical distribution** with $|\Sigma|$ categories, i.e., a vector of probabilities whose components correspond to the probabilities of individual categories. Perhaps the simplest way to represent a categorical distribution is as a vector on a probability simplex.

Definition 3.1.8 (Probability Simplex). *A **probability simplex** Δ^{D-1} is the set of non-negative vectors \mathbb{R}^D whose components sum to 1:*

$$\Delta^{D-1} \stackrel{\text{def}}{=} \left\{ \mathbf{x} \in \mathbb{R}^D \mid x_d \geq 0, d = 1, \dots, D \text{ and } \sum_{d=1}^D x_d = 1 \right\} \quad (3.23)$$

So far, we have framed p_{SM} as a function assigning the conditional distribution over $\bar{\mathbf{y}}$ to each string $\bar{\mathbf{y}}$. The definition of a simplex means that we can more formally express p_{SM} as a **projection** from the Hilbert space of the context representations to $\Delta^{|\Sigma|-1}$, i.e., $p_{\text{SM}}: \mathbb{V} \rightarrow \Delta^{|\Sigma|-1}$. Yet all we have discussed so far is creating a vector $\mathbf{E} \text{enc}(\bar{\mathbf{y}})$ that encodes symbol–context similarities— $\mathbf{E} \text{enc}(\bar{\mathbf{y}})$ is not necessarily on the probability simplex $\Delta^{|\Sigma|-1}$. To address this issue, we turn to projection functions:

Definition 3.1.9 (Projection Function). *A **projection function** $\mathbf{f}_{\Delta^{D-1}}$ is a mapping from a real-valued Hilbert space \mathbb{R}^D to the probability simplex Δ^{D-1}*

$$\mathbf{f}_{\Delta^{D-1}}: \mathbb{R}^D \rightarrow \Delta^{D-1}, \quad (3.24)$$

which allows us to define a probability distribution according to $\mathbf{E} \text{enc}(\bar{\mathbf{y}})$:

$$p_{\text{SM}}(\bar{\mathbf{y}} \mid \bar{\mathbf{y}}) = \mathbf{f}_{\Delta^{|\Sigma|-1}}(\mathbf{E} \text{enc}(\bar{\mathbf{y}}))_{\bar{\mathbf{y}}} \quad (3.25)$$

Clearly, we still want the projection of $\mathbf{E} \text{enc}(\bar{\mathbf{y}})$ onto $\Delta^{|\Sigma|-1}$ to maintain several attributes of the original vector—otherwise, we will lose the notion of compatibility that $\mathbf{E} \text{enc}(\bar{\mathbf{y}})$ inherently encodes. However, $\mathbf{f}_{\Delta^{|\Sigma|-1}}$ must satisfy several additional criteria in order to map onto a valid point in $\Delta^{|\Sigma|-1}$. For example, the inner product of two vectors (and consequently $\mathbf{E} \text{enc}(\bar{\mathbf{y}})$) is not necessarily positive—yet all points in $\Delta^{|\Sigma|-1}$ are positive (see Definition 3.1.8). These characteristics motivate the use of a projection function that is both monotonic and positive everywhere. Thus, one clear choice is to base our projection function on the exponential function, i.e.,

$$\mathbf{f}_{\Delta^{|\Sigma|-1}}(\mathbf{E} \text{enc}(\bar{\mathbf{y}})) \propto \exp(\mathbf{E} \text{enc}(\bar{\mathbf{y}})). \quad (3.26)$$

To make a function of the form in Eq. (3.26) a valid projection function, we now simply have to ensure that the output of $\mathbf{f}_{\Delta^{D-1}}$ sums to 1, which can easily be accomplished by *re-normalizing* the vector of exponentiated values by their sum. This brings us to the main star of this subsection: the softmax.

While we simply motivated its introduction by chasing our goal of ending up on the probability simplex, the origin of the softmax function goes back to the Boltzmann distribution from statistical mechanics introduced in the mid-1800s by Boltzmann (1868). It was then studied intensely and popularized by Gibbs (1902). It was originally introduced as a way to convert the energy function of the Boltzmann distribution into a probability distribution.⁴ Yet now, for reasons we will see in this subsection, the softmax is the predominant choice of projection function in machine learning applications.

Formally, the softmax is often defined in terms of a temperature parameter τ as follows.

Definition 3.1.10 (Softmax). *Let $\tau \in \mathbb{R}_+$ be the **temperature**. The **softmax** at temperature τ is the projection function defined as:*

$$\text{softmax}(\mathbf{x})_d \stackrel{\text{def}}{=} \frac{\exp\left[\frac{1}{\tau}x_d\right]}{\sum_{j=1}^D \exp\left[\frac{1}{\tau}x_j\right]}, \text{ for } d = 1, \dots, D \quad (3.27)$$

where the temperature parameter τ gives us a mechanism for controlling the entropy of the softmax function by *scaling* the individual scores in the input vector before their exponentiation. In the context of the Boltzmann distribution, it was used to control the “randomness” of the system: When the temperature is high, the softmax function outputs a more uniform probability distribution whose probabilities are relatively evenly spread out among the different categories. When the temperature is low, the softmax function outputs a peaked probability distribution, where the probability mass is concentrated on the most likely category. In the limit, as we take τ to the edge of the possible values it can assume, the following properties hold:

Theorem 3.1.2.

$$\lim_{\tau \rightarrow \infty} \text{softmax}(\mathbf{x}) = \frac{1}{D} \mathbf{1} \quad (3.28)$$

$$\lim_{\tau \rightarrow 0^+} \text{softmax}(\mathbf{x}) = \mathbf{e}_{\text{argmax}(\mathbf{x})}, \quad (3.29)$$

where \mathbf{e}_d denotes the d^{th} basis vector in \mathbb{R}^D , $\mathbf{1} \in \mathbb{R}^D$ the vector of all ones, and

$$\text{argmax}(\mathbf{x}) \stackrel{\text{def}}{=} \min \left\{ d \mid x_d = \max_{d=1, \dots, D} (x_d) \right\}, \quad (3.30)$$

i.e., the index of the maximum element of the vector \mathbf{x} (with the ties broken by choosing the lowest such index). In words, this means that the output of the softmax approaches the uniform distribution as $\tau \rightarrow \infty$ and towards a single mode as $\tau \rightarrow 0^+$.⁵

Proof. Let us first consider the case of $\tau \rightarrow 0^+$. Without loss of generality, let us consider a 2-dimensional vector $\mathbf{x} = [x_1, x_2]^\top$

$$\lim_{\tau \rightarrow 0^+} \text{softmax}(\mathbf{x})_1 = \lim_{\tau \rightarrow 0^+} \frac{\exp\left(\frac{x_1}{\tau}\right)}{\exp\left(\frac{x_1}{\tau}\right) + \exp\left(\frac{x_2}{\tau}\right)} \quad (3.31)$$

$$= \lim_{\tau \rightarrow 0^+} \frac{\exp\left(\frac{x_1}{\tau}\right) \exp\left(-\frac{x_1}{\tau}\right)}{\left(\exp\left(\frac{x_1}{\tau}\right) + \exp\left(\frac{x_2}{\tau}\right)\right) \exp\left(-\frac{x_1}{\tau}\right)} \quad (3.32)$$

$$= \lim_{\tau \rightarrow 0^+} \frac{1}{1 + \exp\left(\frac{x_2 - x_1}{\tau}\right)} \quad (3.33)$$

which leads us to the following definition for element-wise values:

$$\lim_{\tau \rightarrow 0^+} \exp\left(\frac{x_2 - x_1}{\tau}\right) = \begin{cases} 0, & \text{if } x_1 > x_2 \\ 1, & \text{if } x_1 = x_2 \\ \infty, & \text{o.w.} \end{cases} \quad (3.34)$$

⁴This is precisely the connection we mentioned in Definition 2.4.1.

⁵ $\tau \rightarrow 0^+$ denotes the limit from above.

Then the limit of softmax as $\tau \rightarrow 0^+$ is given as

$$\lim_{\tau \rightarrow 0^+} \text{softmax}(\mathbf{x}) = \begin{cases} [1, 0]^\top, & \text{if } x_1 > x_2 \\ [\frac{1}{2}, \frac{1}{2}]^\top, & \text{if } x_1 = x_2 \\ [0, 1]^\top, & \text{o.w.} \end{cases} \quad (3.35)$$

which is equivalent to the argmax operator over \mathbf{x} . The proof extends to arbitrary D -dimensional vectors.

The case of $\tau \rightarrow \infty$ follows similar logic, albeit $\lim_{\tau \rightarrow \infty} \exp\left(\frac{x_2 - x_1}{\tau}\right) = 1$ in all cases. Hence, we get $\lim_{\tau \rightarrow \infty} \text{softmax}(\mathbf{x}) = \frac{1}{D} \mathbf{1}$. ■

The second property, specifically, shows that the softmax function resembles the argmax function as the temperature approaches 0—in that sense, a more sensible name for the function would have been “soft argmax”. We will most often simply take τ to be 1. However, different values of the parameter are especially useful when *sampling* or generating text from the model, as we discuss subsequently.

The output of the softmax is equivalent to the solution to a particular optimization problem, giving it a variational interpretation.

Theorem 3.1.3. *Given a set of real-valued scores \mathbf{x} , the following equality holds*

$$\text{softmax}(\mathbf{x}) = \underset{\mathbf{p} \in \Delta^{D-1}}{\text{argmax}} \left(\mathbf{p}^\top \mathbf{x} - \tau \sum_{d=1}^D p_d \log p_d \right) \quad (3.36)$$

$$= \underset{\mathbf{p} \in \Delta^{D-1}}{\text{argmax}} (\mathbf{p}^\top \mathbf{x} + \tau H(\mathbf{p})) \quad (3.37)$$

This tells us that softmax can be given a variational characterization, i.e., it can be viewed as the solution to an optimization problem.

Proof. Eq. (3.36) can equivalently be written as

$$\text{softmax}(\mathbf{x}) = \underset{\left(\mathbf{p}^\top \mathbf{x} - \tau \sum_{d=1}^D p_d \log p_d \right)}{\text{argmax}} \quad (3.38)$$

$$\text{s.t. } \sum_d p_d = 1 \quad (3.39)$$

from which we can clearly see that the Lagrangian of this optimization problem is $\Lambda = \mathbf{p}^\top \mathbf{x} - \tau \sum_{d=1}^D p_d \log p_d + \lambda \sum_d p_d$. Taking the derivative of Λ with respect to p_d , we see that the optimum of is reached when

$$\frac{\partial \Lambda}{\partial p_d} = v_d - \tau(\log p_d + 1) + \lambda = 0 \quad (3.40)$$

Solving for p_d gives us $p_d = Z \exp(\frac{x_d}{\tau})$, where Z is the normalizing constant that ensures $\sum_d p_d = 1$. This solution is equivalent to performing the softmax operation over \mathbf{x} , as desired. ■

Theorem 3.1.3 reveals an interpretation of the softmax as the projection $\mathbf{p} \in \Delta^{D-1}$ that has the maximal similarity with \mathbf{x} while being regularized to produce a solution with high entropy. Further, from both Definition 3.1.10 and Eq. (3.36), we can see that softmax leads to non-sparse solutions as an entry $\text{softmax}(\mathbf{x})_i$ can only be 0 if $x_d = -\infty$.

In summary, the softmax has a number of desirable properties for use in machine learning settings.

Theorem 3.1.4. *The softmax function with temperature parameter τ exhibits the following properties.*

1. *In the limit as $\tau \rightarrow 0^+$ and $\tau \rightarrow \infty$, the softmax recovers the argmax operator and projection to the center of the probability simplex (at which lies the uniform distribution), respectively.*
2. *$\text{softmax}(\mathbf{x} + c\mathbf{1}) = \text{softmax}(\mathbf{x})$ for $c \in \mathbb{R}$, i.e., the softmax is invariant to adding the same constant to all coordinates in \mathbf{x} .*

3. The derivative of the softmax is continuous and differentiable everywhere; the value of its derivative can be explicitly computed.
4. For all temperatures $\tau \in \mathbb{R}_+$, if $x_i \leq x_j$, then $\text{softmax}(\mathbf{x})_i \leq \text{softmax}(\mathbf{x})_j$. In words, the softmax maintains the rank of \mathbf{x} .

Proof. Property 1. is simply a restatement of Theorem 3.1.2. The proof for property 2. can be shown using simple algebraic manipulation:

$$\text{softmax}(\mathbf{x} + c\mathbf{1})_d = \frac{\exp[\frac{1}{\tau}x_d + c]}{\sum_{j=1}^D \exp[\frac{1}{\tau}x_j + c]} = \frac{\exp[\frac{1}{\tau}x_d] \cdot \exp c}{\sum_{j=1}^D \exp[\frac{1}{\tau}x_j] \cdot \exp c} = \text{softmax}(\mathbf{x})_d \quad (3.41)$$

The derivative of the softmax at position i with respect to the variable at position j is given by

$$\frac{\partial \text{softmax}(\mathbf{x})_i}{\partial x_j} = \frac{\delta_i(j) \cdot \exp(x_i) \sum_k \exp(x_k) - \exp(x_i) \cdot \exp(x_j)}{(\sum_k \exp(x_k))^2} \quad (3.42)$$

where $\delta_i(j)$ is the Dirac Delta function, defined as $\delta_i(j) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else} \end{cases}$. Clearly, Eq. (3.42) is continuous.

Further, it takes on values for all $\mathbf{x} \in \mathbb{R}^d$. Lastly, property 4. follows from the monotonicity of the exp function. ■

There are many other valid projection functions that one could choose from. For example, [Martins and Astudillo \(2016\)](#) introduce the **sparsemax**, which *can* output sparse distributions:

$$\text{sparsemax}(\mathbf{x}) \stackrel{\text{def}}{=} \underset{\mathbf{p} \in \Delta^{D-1}}{\text{argmin}} \|\mathbf{p} - \mathbf{x}\|_2^2 \quad (3.43)$$

In words, sparsemax directly maps \mathbf{x} onto the probability simplex, which often leads to solutions on the boundary, i.e., where at least one entry of \mathbf{p} is 0. [Martins and Astudillo \(2016\)](#) provide a method for computing the closed form solution of this optimization problem in Alg. 1 of their work. [Blondel et al. \(2019\)](#) later introduced a framework that encompasses many different projection functions, which they term regularized prediction functions. Essentially, this framework considers the subset of projection functions that can be written as:

$$\mathbf{f}_{\Delta^{|\Sigma|-1}}(\mathbf{x}) \stackrel{\text{def}}{=} \underset{\mathbf{p} \in \Delta^{D-1}}{\text{argmax}} (\mathbf{p}^\top \mathbf{x} - \Omega(\mathbf{p})) \quad (3.44)$$

where $\Omega: \mathbb{R}^D \rightarrow \mathbb{R}$ is regularization term. For certain choices of Ω , there are straightforward closed-form solutions to Eq. (3.44). For example, as we can see from Eq. (3.36), Eq. (3.44) is equivalent to the softmax when $\Omega(\mathbf{p}) = -\mathbf{H}(\mathbf{p})$, meaning we can compute its closed form using Eq. (3.27). Further, we recover the sparsemax when $\Omega(\mathbf{p}) = -\|\mathbf{p}\|_2^2$, which likewise has a closed-form solution. The notion of regularizing \mathbf{p} may be unintuitive at first, but we can view it as trying to balance out the “suitability” term $\mathbf{p}^\top \mathbf{x}$ with a “confidence” term $\Omega(\mathbf{p})$, which should be smaller when \mathbf{p} is “uncertain.” We point the interested reader to the comprehensive work of [Blondel et al. \(2019\)](#) for further elaboration.

So why aren’t these other projection functions more widely employed in machine learning frameworks? First, not all choices of Ω lead to closed-form solutions; further, not all meet the desirable criterion listed in Theorem 3.1.4. For example, the sparsemax is not everywhere differentiable, meaning that one could not simply use out-of-the-box automatic differentiation frameworks when training a model using the sparsemax as its projection function. Rather one would have to specify its gradient explicitly.

Theorem 3.1.5. *The derivative of the the sparsemax with respect to its input \mathbf{x} is as follows:*

$$\frac{\partial \text{sparsemax}(\mathbf{x})_i}{\partial x_j} = \begin{cases} \delta_{ij} - \frac{1}{S(\mathbf{x})} & \text{if } i, j \in S(\mathbf{x}) \\ 0 & \text{else} \end{cases} \quad (3.45)$$

Proof. See [Martins and Astudillo \(2016\)](#). ■

To conclude, projection functions, together with symbol representations and the representation function enc , give us the tools to define a probability distribution over next symbols that encodes complex linguistic interactions. We now bring all the components together into the locally normalized modeling framework in the next section.

3.1.4 Representation-based Locally Normalized Models

With these tools at hand, we now define representation-based locally normalized language models.

Definition 3.1.11 (Representation-Based Locally Normalized Model). *Let enc be an encoding function. A representation-based locally normalized model is a model of the following form:*

$$p_{SM}(\bar{y}_t \mid \bar{y}_{<t}) \stackrel{\text{def}}{=} \mathbf{f}_{\Delta|\Sigma|-1}(\mathbf{E} \text{enc}(\bar{y}_{<t}))_{\bar{y}_t} \quad (3.46)$$

where unless otherwise stated, we assume $\mathbf{f}_{\Delta|\Sigma|-1} = \text{softmax}$. It defines the probability of an entire string $\mathbf{y} \in \Sigma^*$ as

$$p_{LN}(\mathbf{y}) \stackrel{\text{def}}{=} p_{SM}(\text{EOS} \mid \mathbf{y}) \prod_{t=1}^T p_{SM}(y_t \mid \mathbf{y}_{<t}) \quad (3.47)$$

where $y_0 \stackrel{\text{def}}{=} \text{BOS}$.

Alternatively, we could also include an additive **bias** term \mathbf{b} as part of the projection function $\mathbf{f}_{\Delta|\Sigma|-1}$ in the definition of the conditional distribution $p_{SM}(\bar{y}_t \mid \bar{y}_{<t})$, i.e., $p_{SM}(\bar{y}_t \mid \bar{y}_{<t}) = \mathbf{f}_{\Delta|\Sigma|-1}(\mathbf{E} \text{enc}(\bar{y}_{<t}) + \mathbf{b})_{\bar{y}_t}$. However, note that the bias term can be absorbed into the encoding function enc , meaning that we can assume the form Eq. (3.46) without loss of generality. In representation-based language models, $\mathbf{e}(\bar{y})$ and $\text{enc}(\mathbf{y})$ carry all the necessary information to determine how probable individual symbols y are given the context \mathbf{y} . Therefore, the design choices of $\mathbf{e}(\bar{y})$ and $\text{enc}(\mathbf{y})$ are crucial when building language models this way. Indeed, a large portion of the discussion in the remainder of the notes will center around how to build good representations of the context and individual symbols.

3.1.5 Tightness of Softmax Representation-based Models

Having introduced representation-based language models, we can now state a very general result about the tightness of such models. It connects the notion of tightness to the intuition about the “compatibility” of symbols to the context—namely, the compatibility of the EOS symbol to the context (compared to the compatibility of all other symbols). The compatibility is here captured by the distance of the representation of the EOS symbol to the representation of the other symbols—if this distance grows slowly enough with respect to t (modulo the norm of the context representation), the model is tight.

Theorem 3.1.6 (Proposition 5.9 in Du et al., 2022). *Let p_{SM} be a representation-based sequence model over the alphabet Σ , as defined in Definition 3.1.11. Let*

$$s \stackrel{\text{def}}{=} \sup_{y \in \Sigma} \|\mathbf{e}(y) - \mathbf{e}(\text{EOS})\|_2, \quad (3.48)$$

i.e., the largest distance to the representation of the EOS symbol, and

$$z_{max} \stackrel{\text{def}}{=} \max_{\mathbf{y} \in \Sigma^t} \|\text{enc}(\mathbf{y})\|_2, \quad (3.49)$$

i.e., the maximum attainable context representation norm for contexts of length t . Then the locally normalized model p_{LN} induced by p_{SM} is tight if

$$sz_{max} \leq \log t. \quad (3.50)$$

Proof. Let $\mathbf{x}_t(\boldsymbol{\omega})$ be the random variable that is equal to the t^{th} token in an outcome $\boldsymbol{\omega} \in \Omega$. Then for an arbitrary $t \in \mathbb{N}$ and any $\mathbf{y} \in \Sigma^t$, we have:

$$\mathbb{P}(\mathbf{x}_t = \text{EOS} \mid \mathbf{x}_{<t} = \mathbf{y}) = \frac{\exp[\mathbf{e}(\text{EOS})^\top \text{enc}(\mathbf{y})]}{\sum_{\mathbf{y} \in \Sigma} \exp[\mathbf{e}(\mathbf{y})^\top \text{enc}(\mathbf{y})]} \quad (3.51a)$$

$$= \frac{1}{\frac{\sum_{\mathbf{y} \in \Sigma} \exp[\mathbf{e}(\mathbf{y})^\top \text{enc}(\mathbf{y})]}{\exp[\mathbf{e}(\text{EOS})^\top \text{enc}(\mathbf{y})]}} \quad (3.51b)$$

$$= \frac{1}{1 + \sum_{\mathbf{y} \in \Sigma} \exp[(\mathbf{e}(\mathbf{y}) - \mathbf{e}(\text{EOS}))^\top \text{enc}(\mathbf{y})]} \quad (3.51c)$$

$$\geq \frac{1}{1 + \sum_{\mathbf{y} \in \Sigma} \exp[\|\mathbf{e}(\mathbf{y}) - \mathbf{e}(\text{EOS})\|_2 \|\text{enc}(\mathbf{y})\|_2]} \quad (\text{Cauchy-Schwarz}) \quad (3.51d)$$

$$\geq \frac{1}{1 + \sum_{\mathbf{y} \in \Sigma} \exp[k \|\text{enc}(\mathbf{y})\|_2]} \quad (3.51e)$$

$$= \frac{1}{1 + |\Sigma| \exp[\|\text{enc}(\mathbf{y})\|_2]} \quad (3.51f)$$

Now define $z_{\max} \stackrel{\text{def}}{=} \sup_{\mathbf{y} \in \Sigma^t} \|\text{enc}(\mathbf{y})\|_2$. We then have that $\forall t \in \mathbb{N}$ and $\forall \mathbf{y} \in \Sigma^t$:

$$\mathbb{P}(\mathbf{x}_t = \text{EOS} \mid \mathbf{x}_{<t} = \mathbf{y}) \geq \frac{1}{1 + |\Sigma| \exp(k z_{\max})} \quad (3.52)$$

Now, by Proposition 2.5.6, we have that if $\sum_{t=1}^{\infty} \frac{1}{1 + |\Sigma| \exp(k z_{\max})}$ diverges, then the language model is tight. We will show that if we have that $\exists N \in \mathbb{N}$ such that $\forall t \geq N$, $k z_{\max} \leq \log t$, then the sequence model must be tight.

First, note that $\lim_{t \rightarrow \infty} \frac{1}{t} \frac{1 + |\Sigma| t}{1} = \lim_{t \rightarrow \infty} \frac{1}{t} + |\Sigma| = |\Sigma| \in (0, \infty)$. Hence, by the limit comparison test, since $\sum_{t=1}^{\infty} \frac{1}{t}$ diverges, this means $\sum_{t=1}^{\infty} \frac{1}{1 + |\Sigma| t}$ must also diverge.

Now, suppose that $k z_{\max} \leq \log t$ for all $t \geq N$. This implies that for $t \geq N$ we have $\frac{1}{1 + |\Sigma| \exp(k z_{\max})} \geq \frac{1}{1 + |\Sigma| t}$, which combined with the above and the comparison test, implies that $\sum_{t=N}^{\infty} \frac{1}{1 + |\Sigma| \exp(k z_{\max})}$ diverges. This in turn means that $\sum_{t=1}^{\infty} \frac{1}{1 + |\Sigma| \exp(k z_{\max})}$ diverges. Hence, if $k z_{\max} \leq \log t$ for all $t \geq N$ for some $N \in \mathbb{N}$, then the language model is tight. ■

Theorem 3.1.6 is a generalization of the following result from Welleck et al. (2020).

Theorem 3.1.7. *A locally-normalized representation-based language model, as defined in Definition 3.1.11, with uniformly bounded $\|\text{enc}(\mathbf{y})\|_p$ (for some $p \geq 1$) is tight.*

For most of the language models that we consider, $\text{enc}(\mathbf{y})$ is bound due to the choice of activation functions. In turn, $\mathbf{E} \text{enc}(\bar{\mathbf{y}}_{<t})$ is bounded for all $\bar{\mathbf{y}}$. Further, by the definition of the softmax, $\mathbf{f}_{\Delta[\Sigma]-1}(\mathbf{E} \text{enc}(\bar{\mathbf{y}}_{<t}))_{\text{EOS}} > \eta$ for some constant η .

This concludes our investigation of general representation-based models. The next section discusses learning parametrized models (as a special case, also symbol and context representations).

3.2 Estimating a Language Model from Data

The **language modeling task** refers to any attempt to estimate the parameters⁶ of a model p_M of the ground-truth probability distribution over natural language strings p_{LM} using data $\mathcal{D} = \{\mathbf{y}^{(n)}\}_{n=1}^N$, where we assume samples $\mathbf{y}^{(n)}$ were generated according to p_{LM} . This task is often treated as an optimization problem. Here we will discuss the various components of this optimization problem, primarily the objective and the algorithm used to perform optimization. Note that the material covered here corresponds to what is colloquially referred to as pre-training. The learning paradigm for fine-tuning a language model for a downstream task will be covered later in the course.

3.2.1 Data

In this course, we consider objectives that are defined in terms of data \mathcal{D} . Therefore, we will first discuss the nature of this data which, more precisely, is a corpus of texts. Following the notation used throughout the rest of these notes, let Σ be an alphabet. A **corpus** $\mathcal{D} = \{\mathbf{y}^{(n)}\}_{n=1}^N \subset \Sigma^*$ is a collection of N strings. We will use the terms corpus and dataset interchangeably throughout this section. We make the following assumption about the data-generating process of \mathcal{D} :

Assumption 3.2.1. *The strings $\mathbf{y}^{(n)}$ in our corpus \mathcal{D} are generated independently and identically distributed (i.i.d.) by some unknown distribution p_{LM} .*

Note that $\mathbf{y}^{(n)}$ are strings of an arbitrary length; they can be single words, sentences, paragraphs, or even entire documents depending on how we choose Σ . For example, often our models' architectural designs make them unable to process document-length strings efficiently, e.g., they might not fit into a context window that can be reasonably processed by a transformer language model; we will elaborate on this statement in our discussion of transformers in §5.2. Thus in practice, we often chunk documents into paragraphs that we treat as separate data points.⁷ This means that our model may not be able to learn properties of language such as discourse structure.

3.2.2 Language Modeling Objectives

Similarly to many other machine learning tasks, we can cast our problem as the *search* for the best model p_M of the ground-truth distribution over strings p_{LM} . In order to make this search tractable, we must limit the models p_M that we consider. Explicitly, we make the following assumption:

Assumption 3.2.2. *p_{LM} is a member of the parameterized family of models $\{p_\theta \mid \theta \in \Theta\}$, the set of all distributions representable by parameters θ in a given parameter space Θ .*

As concrete examples, θ could be the conditional probabilities in a simple, standard n -gram model for a given prefix of size $n - 1$, i.e., θ is $n - 1$ simplices of size $|\Sigma|$.⁸ As another example, θ could be the weights of a neural network; the set Θ would then cover all possible valid weight matrices that could parameterize our model.

Assumption 3.2.2 implies that we can equivalently write p_{LM} as p_{θ^*} for certain (unknown) parameters $\theta^* \in \Theta$.⁹ Further, an arbitrary model p_M from this hypothesis space with parameters θ can be written as p_θ ; we will use this notation for the remainder of the chapter to make the parameterization of our distribution explicit. We now turn to the general framework for choosing the best parameters $\theta \in \Theta$ so that our model p_θ serves as a good approximation of p_{θ^*} .¹⁰

⁶Most of this course focuses on the parametric case, i.e., where p_M is governed by a set of parameters θ . However, we will briefly touch upon various non-parametric language models.

⁷This practice technically breaks Assumption 3.2.1, yet the negative (empirically-observed) effects of this violation are minimal and perhaps outweighed by the additional data it allows us to make use of.

⁸One might be tempted to assume we only need $|\Sigma| - 1$ parameters per simplex, but we condition over $\bar{\Sigma}$ classes per prefix position.

⁹We discuss the implications of the case that $p_{LM} \notin \{p_\theta \mid \theta \in \Theta\}$ later in this section.

¹⁰The modeling paradigms that we will discuss in this section are predominantly generative, i.e., these models try to learn the underlying distribution of the data rather than the boundaries between different classes or categories. The implication is that parameter estimation in language modeling typically makes use of *unannotated* text data, and is therefore sometimes referred to as **self-supervised**.

General Framework

We search for model parameters $\hat{\theta} \in \Theta$ such that the model induced by those parameters maximizes a chosen objective, or alternatively, minimizes some loss function $\ell : \Theta \times \Theta \rightarrow \mathbb{R}_{\geq 0}$. This loss can be used to measure the quality of this model as an approximation to p_{θ^*} . In simple math, we search for the solution.

$$\hat{\theta} \stackrel{\text{def}}{=} \underset{\theta \in \Theta}{\operatorname{argmin}} \ell(\theta^*, \theta) \quad (3.53)$$

where our loss function is chosen with the following principle in mind

Principle 3.2.1 (Proximity Principle). *We seek a model p_{θ} that is “close” to p_{θ^*} .*

That is, we choose our loss function to be a measure M of the difference between a distribution parameterized by θ and one parameterized by the true θ^* , i.e., those of our ground-truth distribution. Yet we are immediately faced with a problem: computing an arbitrary M between θ and θ^* (or at least the distributions induced by these sets of parameters) requires knowledge of both, the latter for which we only have samples $\mathbf{y}^{(n)} \in \mathcal{D}$. We will therefore use our corpus \mathcal{D} as an approximation to p_{θ^*} , which is typically implemented by representing \mathcal{D} as an empirical distribution—a collection of Dirac Delta functions—which we will denote as \widetilde{p}_{θ^*} . Formally, we define

$$\widetilde{p}_{\theta^*}(\mathbf{y}) \stackrel{\text{def}}{=} \frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{y}^{(n)}}(\mathbf{y}) \quad (3.54)$$

where the Dirac Delta function $\delta_{x'}(x) = \begin{cases} 1 & \text{if } x = x' \\ 0 & \text{else} \end{cases}$ is essentially a point mass with all probability on x' .

We can decompose this definition over symbols in our strings as well. I.e., we can compute

$$\widetilde{p}_{\theta^*}(y_t | \mathbf{y}_{<t}) = \frac{1}{N_{\mathbf{y}_{<t}}} \sum_{n=1}^N \delta_{y_t^{(n)} | \mathbf{y}_{<t}^{(n)}}(y_t | \mathbf{y}_{<t}) \quad (3.55)$$

where $N_{\mathbf{y}_{<t}} \stackrel{\text{def}}{=} \sum_{n=1}^N \mathbb{1}\{\mathbf{y}_{<t}^{(n)} = \mathbf{y}_{<t}\}$. Note that we can likewise define Eq. (3.55) in terms of the one-hot encodings of symbols, i.e., using the definition in Example 3.1.2: $\widetilde{p}_{\theta^*}(\cdot | \mathbf{y}_{<t}) = \frac{1}{N_{\mathbf{y}_{<t}}} \sum_{n=1}^N \llbracket y_t^{(n)} \rrbracket \mathbb{1}\{\mathbf{y}_{<t}^{(n)} = \mathbf{y}_{<t}\}$. In fact, the empirical distribution is often also referred to in machine learning as the one-hot encoding of a dataset.

Now that we are equipped with methods for representing both p_{θ^*} and p_{θ} , we can define a loss function for approximating p_{θ^*} using p_{θ} .

Cross-Entropy. A natural choice for a loss function is cross-entropy, a measure of the difference between two probability distributions, which has its roots in information theory (Shannon, 1948). Specifically, in Eq. (3.53), we take $\ell(\theta^*, \theta) = H(\widetilde{p}_{\theta^*}, p_{\theta})$ where the definition of the cross-entropy H between distributions p_1 (with support \mathcal{Y}) and p_2 is as follows:

$$H(p_1, p_2) = - \sum_{\mathbf{y} \in \mathcal{Y}} p_1(\mathbf{y}) \log p_2(\mathbf{y}) \quad (3.56)$$

Further, most of the models that we will encounter in this course are locally normalized. Thus, it is more common to see cross-entropy expressed as

$$H(p_1, p_2) = - \sum_{\mathbf{y} \in \mathcal{Y}} \sum_{t=1}^T p_1(y_t^{(n)}) \log p_2(y_t | \mathbf{y}_{<t}). \quad (3.57)$$

Note that cross-entropy is not symmetric, i.e., $H(p_1, p_2) \neq H(p_2, p_1)$. To motivate cross-entropy as a loss function, as well as the intuitive difference between the two argument orderings, we turn to coding theory, a sub-field of information theory. In words, the cross-entropy between two probability distributions is the expected number of bits needed to encode an event $\mathbf{y} \in \mathcal{Y}$ from p_1 when using the optimal encoding scheme corresponding to distribution p_2 . Importantly, the optimal encoding scheme for p_1 uses $\log p_1(\mathbf{y})$ bits to encode an event \mathbf{y} that occurs with probability $p_1(\mathbf{y})$, implying that the minimal cross-entropy is achieved when $p_1 = p_2$. This characteristic of cross-entropy motivates another metric: the KL divergence D_{KL} .

KL Divergence. A **divergence measure** is a measure of statistical distance¹¹ between two probability distributions. The KL divergence is defined as:

$$D_{\text{KL}}(p_1 \parallel p_2) = \sum_{\mathbf{y} \in \mathcal{Y}} p_1(\mathbf{y}) \log p_2(\mathbf{y}) - p_1(\mathbf{y}) \log p_1(\mathbf{y}) \quad (3.58)$$

The KL divergence can intuitively be viewed as the cross-entropy shifted by the expected number of bits used by the optimal encoding scheme for p_1 , i.e., it is the *additional* number of expected bits needed to encode events from p_1 when using our encoding scheme from p_2 . Indeed, taking $\ell(\boldsymbol{\theta}^*, \boldsymbol{\theta}) = D_{\text{KL}}(\widetilde{p}_{\boldsymbol{\theta}^*} \parallel p_{\boldsymbol{\theta}})$ should lead to the same solution as taking $\ell(\boldsymbol{\theta}^*, \boldsymbol{\theta}) = H(\widetilde{p}_{\boldsymbol{\theta}^*}, p_{\boldsymbol{\theta}})$ because the $\widetilde{p}_{\boldsymbol{\theta}^*}(\mathbf{y}) \log \widetilde{p}_{\boldsymbol{\theta}^*}(\mathbf{y})$ term is constant with respect to model parameters $\boldsymbol{\theta}$.

Relationship to Maximum Likelihood Estimation

An alternative way that we could frame our search for model parameters $\hat{\boldsymbol{\theta}} \in \boldsymbol{\Theta}$ is in terms of data likelihood. Formally, the **likelihood** of the corpus \mathcal{D} under the distribution $p_{\boldsymbol{\theta}}$ is the joint probability of all $\mathbf{y}^{(n)}$:

$$L(\boldsymbol{\theta}) = \prod_{n=1}^N p_{\boldsymbol{\theta}}(\mathbf{y}^{(n)}). \quad (3.59)$$

The principle of maximum likelihood then dictates:

Principle 3.2.2 (Maximum Likelihood). *The optimal parameters for a model are those that maximize the likelihood of observing the given data under that model. Formally:*

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta} \in \boldsymbol{\Theta}}{\operatorname{argmax}} \mathcal{L}(\boldsymbol{\theta}) \quad (3.60)$$

Note that in practice, we typically work with the **log-likelihood** $\mathcal{L}(\boldsymbol{\theta}) = \log L(\boldsymbol{\theta})$ rather than the likelihood for a number of reasons, e.g., it is convex and more numerically stable given the small probabilities we encounter when using L and the finite precision of the computing frameworks that we employ. Since log is a monotonically increasing function, this would not change the solution to Eq. (3.60). Further, as is the case with Eq. (3.57), we decompose our loss over symbol-level distributions.

Notably, in our setting, finding parameters that maximize data log-likelihood is equivalent to finding those that minimize cross-entropy. We show this equivalence below.

Proposition 3.2.1. *The optimal parameters under Eq. (3.60) are equivalent to the optimal parameters when solving for Eq. (3.53) with the cross-entropy loss between the empirical distribution $\widetilde{p}_{\boldsymbol{\theta}^*}$ and the model $p_{\boldsymbol{\theta}}$.*

Proof. Under the standard practice of taking $0 \log(0) = 0$, the only elements of \mathcal{Y} that make a nonzero contribution to $H(\widetilde{p}_{\boldsymbol{\theta}^*}, p_{\boldsymbol{\theta}})$ are sequences in the support of $\widetilde{p}_{\boldsymbol{\theta}^*}$, making summing over \mathcal{Y} equivalent to summing over \mathcal{D} :

$$H(\widetilde{p}_{\boldsymbol{\theta}^*}, p_{\boldsymbol{\theta}}) = - \sum_{\mathbf{y} \in \Sigma^*} \widetilde{p}_{\boldsymbol{\theta}^*}(\mathbf{y}) \log p_{\boldsymbol{\theta}}(\mathbf{y}) \quad (3.61)$$

$$= - \sum_{\mathbf{y} \in \Sigma^*} \frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{y}^{(n)}}(\mathbf{y}) \log p_{\boldsymbol{\theta}}(\mathbf{y}) \quad (3.62)$$

$$= - \sum_{\mathbf{y} \in \Sigma^*} \frac{1}{N} \mathbb{1}\{\mathbf{y}^{(n)} \in \mathcal{D}\} \log p_{\boldsymbol{\theta}}(\mathbf{y}) \quad (3.63)$$

$$\propto - \sum_{\mathbf{y} \in \mathcal{D}} \log p_{\boldsymbol{\theta}}(\mathbf{y}) \quad (3.64)$$

$$= -\mathcal{L}(\boldsymbol{\theta}) \quad (3.65)$$

Thus, we can see that the objectives are equivalent, up to a multiplicative constant that is independent of model parameters. ■

¹¹Divergences are not technically distances because they are not symmetric, i.e., it may be the case for divergence measure D and probability distributions p and q that $D(p \parallel q) \neq D(q \parallel p)$. However, they do meet the criteria that $D(p \parallel q) \geq 0 \forall p, q$ and $D(p \parallel q) = 0 \iff p = q$.

The equivalence of cross-entropy, D_{KL} divergence, and maximum likelihood as learning objectives provides intuition about our many goals when learning p_θ : (1) we want a close (w.r.t. a given metric) approximation of the data-generating distribution, and (2) this approximation should place high probability on samples of real language data.

Properties of $\hat{\theta}$ under the cross-entropy loss. Assumption 3.2.2 may feel quite strong, as it implies we know a great deal about the nature of p_{LM} . However, it allows us to prove the optimality of p_θ under certain conditions.

Theorem 3.2.1. *Consider that our loss function $\ell(\theta^*, \theta) = H(\tilde{p}_\theta, p_\theta)$ (or equivalently that $\ell(\theta^*, \theta) = D_{\text{KL}}(\tilde{p}_\theta \parallel p_\theta)$). Given Assumption 3.2.1 and that the minimizer p_θ of $H(\tilde{p}_\theta, p_\theta)$ is unique, then under certain (quite strong) regularity conditions on $\{p_\theta \mid \theta \in \Theta\}$, $\hat{\theta}$ is a consistent estimator, i.e., it converges to θ^* in probability as $n \rightarrow \infty$.*

Arguably, in practice, Assumption 3.2.2 does not hold; we often make some incorrect modeling assumptions. Naturally, this raises the following question: If we *misspecify* the family of models that p_{LM} belongs to, i.e., $p_{\text{LM}} \notin \{p_\theta \mid \theta \in \Theta\}$, then is our optimal model p_θ under the cross-entropy loss at all meaningful? Fortunately, the answer here is yes. In this case, we can interpret p_θ as a projection of p_{LM} onto the manifold of parametric models $\{p_\theta \mid \theta \in \Theta\}$. This projection is formally known as an **information projection** (Nielsen, 2018), which while we do not cover formally here, we can intuit as a mapping of p_{LM} onto its “closest” point in $\{p_\theta \mid \theta \in \Theta\}$. In this setting, using different metrics M leads to different definitions of closeness, which in turn means that optimal models under different M exhibit different properties.

Potential drawbacks of cross-entropy loss. A closer inspection of Eq. (3.56) reveals that, when we use $H(\tilde{p}_\theta, p_\theta)$ as our loss function, p_θ must put probability mass on *all* samples $\mathbf{y}^{(n)}$ in the support of \tilde{p}_θ ; otherwise, our loss is infinite. Since the model is not explicitly penalized for extraneous coverage, it will thus resort to placing mass over all of Σ^* to avoid such gaps;¹² this is sometimes referred to as *mean-seeking* behavior. In practice, this means that sequences of symbols that one might qualitatively describe as gibberish are assigned nonzero probability by p_θ . It is unclear whether this is a desirable property under a language model. While perhaps useful when using such a model to assign probabilities to strings—in which case, we might be more interested in how strings’ probabilities rank against each other and may not want to write off any string as completely improbable—it could prove problematic when generating strings from these models, a topic covered later in this course.

Teacher Forcing. The loss functions that we have considered thus far are all based on our model’s predictions conditioned on prior context. Here we are faced with a choice during training: we could either use the model’s predictions from the previous time step(s) $p_\theta(\cdot \mid \hat{\mathbf{y}}_{<t})$ (e.g., the most probable symbols) as the prior context or use the ground-truth prior context from our data $p_\theta(\cdot \mid \mathbf{y}_{<t})$. The latter method is often referred to as **teacher forcing**: Even if our model makes an incorrect prediction at one step of training, we intervene and provide the correct answer for it to make subsequent predictions with.

From a theoretical perspective, training with the cross-entropy loss mandates that we should use the teacher-forcing approach since each conditional distribution is defined with respect to the ground-truth context; this is elucidated, for example, in Eq. (3.57). Yet such meticulous guidance can lead to poor performance in tasks where the model is required to accurately predict an entire sequence of symbols on its own. For example, in language generation, since the model is not exposed to its own generations during training, small errors in predictions can compound, leading to degenerate text. This problem is known as **exposure bias**. Only the other hand, using previous model outputs in order to make subsequent predictions can lead to serious instability during training, especially if implemented from the start of training. Methods for alleviating exposure bias have been proposed with more stable training dynamics, such as scheduled sampling Bengio et al. (2015), which we discuss in §3.2.2.

¹²This behavior can also be (at least partially) attributed to the softmax used to transform model outputs into a probability distribution over symbols. Since the softmax maps to the interior of the probability simplex, no symbol can be assigned a probability of exactly 0.

Alternative Objectives

Masked Language Modeling. So far, our parameter estimation strategies have made use of the decomposition of $p_{\theta}(\mathbf{y})$ into individual symbol probabilities, conditioned on *prior* symbols, i.e., $p_{\theta}(\mathbf{y}) = \prod_{t=1}^T p_{\theta}(y_t | \mathbf{y}_{<t})$. In other words, we do not give a model both sides of a symbol’s context when asking it to estimate the probability distribution over that symbol. While this paradigm might be more realistic when using a language model for tasks such as generation—for which we may want to generate outputs sequentially to mimic human language production—access to both sides of a symbol’s context could be critical when using the model for tasks such as acceptability judgments or classification. This motivates the use of an alternative objective for parameter estimation.

Similarly to the maximum likelihood objective in Eq. (3.59), we can choose model parameters by optimizing for the per-symbol log-likelihood of a dataset \mathcal{D} , albeit in this case, using *both* sides of the symbol’s context:

$$\mathcal{L}_{\text{MLM}}(\theta) = \sum_{n=1}^N \sum_{t=1}^T \log p_{\theta}(y_t^{(n)} | \mathbf{y}_{<t}^{(n)}, \mathbf{y}_{>t}^{(n)}) \quad (3.66)$$

Eq. (3.66) is sometimes referred to as the **pseudo(log)likelihood** (Besag, 1975), since it gives us an approximation of the true log-likelihood, i.e., $\sum_{t=1}^T \log p_{\theta}(y_t | \mathbf{y}_{<t}, \mathbf{y}_{>t}) \approx \log p_{\theta}(\mathbf{y})$. Pseudolikelihood has its origins in thermodynamics, where it was used as an approximate inference technique for parameter estimation in Ising models. In such situations, computing $p_{\theta}(y_t | \mathbf{y}_{\neq t})$ often proved computationally easier than computing the exact set of conditional probabilities whose product equaled the marginal.

Using Eq. (3.66) as a model’s training objective is also motivated by psychological tests of language understanding—specifically, the Cloze (Taylor, 1953) task in psychology, in which the goal is to predict the omitted symbol from a piece of text that constitutes a logical and coherent completion. For example, in the string

Example 3.2.1. *The students [MASK] to learn about language models.*

we predict *want* or *like* with high probability for the [MASK] position. When used as an objective in NLP, estimating the probability distribution over symbols at the masked position is referred to as masked language modeling; BERT (Devlin et al., 2019) is one well known example of a masked language model. In practice, typically only the distributions over symbols at a percentage of randomly-chosen positions in \mathcal{D} are estimated during training. As mentioned in §2.5, a model whose parameters are estimated with the masked language modeling objective is not a valid language model in the sense of Definition 2.3.7 because it does not provide a valid distribution over Σ^* . Yet, masked language models have become increasingly popular as base models for fine-tuning on certain downstream tasks, where they sometimes lead to superior performance over standard language models.

Other Divergence Measures. From a given hypothesis space (see Assumption 3.2.2), the distribution that minimizes a given divergence measure with p_{θ^*} exhibits certain properties with respect to how probability mass is spread over the support of that distribution. For example, the model p_{θ} that minimizes $D_{\text{KL}}(p_{\theta^*} || p_{\theta})$ exhibits *mean-seeking* behavior, as discussed earlier in this section. These properties have been studied in depth by a number of works (Minka, 2005; Theis et al., 2016; Huszár, 2015; Labeau and Cohen, 2019). The implication of these findings is that, depending on the use case for the model, other divergence measures may be better suited as a learning objective. For example, prior work has noted frequency biases in models estimated using the standard log-likelihood objective, i.e., these models exhibit an inability to accurately represent the tails of probability distributions (Gong et al., 2018). This is particularly relevant in the case of language modeling, as symbol-usage in natural language tends to follow a power-law distribution (Zipf, 1935). Consequently, when we care particularly about accurately estimating the probability of rare words, we may wish to instead use a loss function that prioritizes good estimation of probability distribution tails. On the other hand, in the case of language generation, we may desire models that only assign probability mass to outputs that are highly-likely according to p_{θ^*} , even if this means assigning probabilities of 0 to some outcomes possible under p_{θ^*} . In other words, we may want a model with *mode-seeking* behavior, which is characteristic of models trained to minimize $D_{\text{KL}}(p_{\theta} || p_{\theta^*})$. However, there are a number of computational issues with using other divergence measures—such as general power divergences, reverse D_{KL} divergence,

and total variation distance—for training neural probabilistic models over large supports, making them difficult to work with in practice. For example, we can compute a Monte Carlo estimate of the forward D_{KL} divergence simply by using samples from p_{θ^*} , which is exactly what we have in our dataset. However an unbiased estimator of the reverse D_{KL} divergence would require the ability to query p_{θ^*} for probabilities, which we do not have.

Scheduled Sampling and Alternative Target Distributions. Scheduled sampling (Bengio et al., 2015) is an algorithm proposed with the goal of alleviating exposure bias: after an initial period of training using the standard teacher forcing approach, some percentage of the models’ predictions are conditioned on prior model outputs, rather than the ground-truth context. However, under this algorithm, $\hat{\theta}$ does not lead to a consistent estimator of θ^* (Huszár, 2015). Other methods likewise aim to alleviate the discrepancy between settings during parameter estimation and those at inference time by specifying an alternative target distribution, for example, one that ranks “higher-quality” text as more probable than average-quality text. Ultimately, these methods often make use of techniques developed for reinforcement learning, i.e., the REINFORCE algorithm. These methods fall under the category of fine-tuning criterion, which are discussed later in this course.

Auxiliary Prediction Tasks. Certain works jointly optimize for an additional objective when performing parameter estimation. For example, the parameters for BERT were learned using both the masked language modeling objective as well as a task referred to as next sentence prediction, i.e., given two sentences, estimating the probability that the second sentence followed the first in a document. A number of similar auxiliary tasks have subsequently been proposed, such as symbol frequency prediction or sentence ordering (see Aroca-Ouellette and Rudzicz (2020) for summary). However, these tasks do not have a formal relationship to language modeling and it is unclear what their effects are on a model’s ability to serve as a valid probability distribution over strings. They likely lead to models that no longer fulfill the formal criteria of §2.5.4.

3.2.3 Parameter Estimation

Given a loss function ℓ and a parameter space Θ from which to choose model parameters, we are now tasked with *finding* the parameters θ , i.e., solving Eq. (3.53). For the class of models that we consider (those parameterized by large neural networks), finding an exact solution analytically would be impractical, if not impossible. Thus, we must resort to numerical methods, where we find approximately optimal parameters by iterating over solutions. This is known as **parameter estimation**, or more colloquially as **training** our model.

Here we will review the various components of training a language model from start to finish. Many of the techniques used for training language models are generally applicable machine learning techniques, e.g., gradient-descent algorithms. Further, these techniques are constantly evolving and often viewed as trade secrets, meaning that entities building and deploying models may not reveal the combination of components that they employed. Thus, we give a more general overview of the design choices involved in parameter estimation, along with the characteristics common to most components.

Data Splitting

In any machine learning setting, we may overestimate model quality if we evaluate solely on its performance w.r.t. the data on which its parameters were estimated. While we can often construct a model that performs arbitrarily well on a given dataset, our goal is to build a model that generalizes to unseen data. Thus, it is important to measure the final performance of a model on data that has had no influence on the choice of model parameters.

This practice can be accomplished simply by splitting the data into several sets. The two basic data splits are a training set $\mathcal{D}_{\text{train}}$ and test set $\mathcal{D}_{\text{test}}$; as the names imply, the training set is used during parameter estimation while the test set is used for evaluating final performance. When samples from $\mathcal{D}_{\text{test}}$ can be found in $\mathcal{D}_{\text{train}}$, we call this **data leakage**. The training set can be further divided to produce a validation set \mathcal{D}_{val} . Typically, \mathcal{D}_{val} is not used to define the objective for which parameters are optimized. Rather, it serves as a check during training for the generalization abilities of a model, i.e., to see whether the model has started

overfitting to the training data. The validation set can be used, e.g., to determine when to stop updating parameters.

Numerical Optimization

From a starting point $\theta_0 \in \Theta$ chosen according to our initialization strategy, we want to find $\hat{\theta}$ in an efficient manner. This is where numerical **optimization algorithms** come into play—a precise set of rules for choosing how to move within Θ in order to find our next set of parameters. The output of a numerical optimization algorithm is a sequence of iterates $\{\theta_s\}_{s=0}^T$, with the property that as $T \rightarrow \infty$ we find the minimizer of our objective ℓ . Ideally, even after a finite number of iterations, we will be sufficiently close to $\hat{\theta}$.

The basic algorithm for searching the parameter space for $\hat{\theta}$ follows a simple formula: starting from $\theta_0 \in \Theta$, we iteratively compute $\theta_1, \theta_2, \dots$ as

$$\theta_{s+1} = \theta_s + \text{update magnitude} \times \text{update direction}. \quad (3.67)$$

where the update added to θ_s to obtain θ_{s+1} is intended to move us closer to $\hat{\theta}$. Once some maximum number of updates S or a pre-defined desideratum has been met, e.g., our loss has not improved in subsequent iterations, we stop and return the current set of parameters. Many of the numerical optimization techniques in machine learning are gradient-based, i.e., we use the gradient of the objective with respect to current model parameters (denoted as $\nabla_{\theta_s} \ell(\theta_s)$) to determine our update direction. Standard vanilla gradient descent takes the form of §3.2.3, where the learning rate schedule $\eta = \langle \eta_0, \dots, \eta_T \rangle$ determines the step size of our parameter update in the loss minimizing direction—there is an inherent trade-off between the rate of convergence and overshooting—and the stopping criterion C determines whether we can terminate parameter updates before our maximum number of iterations S . In vanilla gradient descent, we set $\eta = c \cdot 1$

Algorithm 1 Gradient descent for parameter optimization.

Input: ℓ objective

θ_0 initial parameters

η learning rate schedule

$C: \ell \times \Theta \times \Theta \rightarrow \{\text{True}, \text{False}\}$ stopping criterion

1. **for** $s = 0, \dots, S$:
 2. $\theta_{s+1} \leftarrow \theta_s - \eta_s \cdot \nabla_{\theta} \ell(\theta_s)$
 3. **if** $C(\ell, \theta_s, \theta_{s-1})$:
 4. **break**
 5. **return** θ_s
-

for some constant c and $C(\ell, \theta_s, \theta_{s-1}) = \mathbb{1}\{|\ell(\theta_s) - \ell(\theta_{s-1})| < \epsilon\}$ for user-chosen ϵ —in words, we stop when the change in loss between parameter updates is below a chosen threshold. In practice, more sophisticated learning rate schedules η , e.g., square-root functions of the timestep (Hoffer et al., 2017) or adaptive functions that take into account model parameter values (Duchi et al., 2011), and stopping criterion C are employed.

Modern training frameworks rely on backpropagation—also known as reverse-mode automatic differentiation (Griewank and Walther, 2008)—to compute gradients efficiently (and, as the name implies, automatically!). In fact, gradients can be computed using backpropagation in the same complexity as evaluation of the original function. We do not provide a formal discussion of backpropagation here but see Griewank and Walther (2008) for this material.

Recall that our loss function—and consequently the gradient of our loss function—is defined with respect to the *entire* dataset. Vanilla gradient descent therefore requires iterating through all of $\mathcal{D}_{\text{train}}$ in order to determine the direction to move parameters U , which is an incredibly time-consuming computation for the large datasets employed in modern machine learning settings. Rather, an optimization algorithm would likely take much less time to converge if it could rapidly compute estimates of the gradient at each step. This is the motivation behind perhaps the most widely employed class of optimization algorithms in machine learning: variations of **stochastic gradient descent** (SGD), such as mini-batch gradient descent. Explicitly, these algorithms make use of the fact that $\mathbb{E}_{\mathcal{D}' \sim \mathcal{D}} \nabla_{\theta} \ell(\theta, \mathcal{D}') = \nabla_{\theta} \ell(\theta, \mathcal{D})$, where in slight abuse of notation, we use $\mathcal{D}' \sim \mathcal{D}$ to signify that the multi-set \mathcal{D}' consists of random i.i.d. samples from \mathcal{D} . Thus we can instead

base our loss ℓ , and consequently U , off of a randomly selected subset of the data.¹³ In practice though, this sample is taken *without* replacement, which breaks the i.i.d. assumption. This in turn implies that our gradient estimates are biased under the mini-batch gradient descent algorithm. However, this bias does not seem to empirically harm the performance of such optimization strategies. Indeed, an entire branch of machine learning called **curriculum learning** focuses on trying to find an optimal data ordering with which to train models to achieve desirable characteristics such as generalization abilities. Even when orderings are randomly selected, the chosen ordering can have a large impact on model performance (Dodge et al. (2020)).

A number of optimization algorithms have since iterated on SGD, e.g., the momentum algorithm (Polyak, 1964). In short, the momentum algorithm computes an exponentially decaying moving average of past gradients, and continues updating parameters in this direction, which can drastically speed up convergence. A widely-employed optimization algorithm called ADAM (Kingma and Ba, 2015) takes a similar approach. Just as in momentum, it computes update directions using a moving average (first moment) of gradients, albeit it additionally makes use of the variance of gradients (second moment) when computing update directions. ADAM is one of the most popular optimization algorithms used for training large language models in modern ML frameworks.

Parameter Initialization

Our search for (approximately) optimal model parameters must start from some point in the parameter space, which we denote as θ_0 . Ideally, starting from any point would lead us to the same solution, or at least to solutions of similar quality. Unfortunately, this is not the case: both training dynamics and the performance of the final model can depend quite heavily on the chosen initialization strategy, and can even have high variance between different runs of the same strategy. This makes sense at some intuitive level though: depending on the learning algorithm, an initial starting point can heavily dictate the amount of searching we will have to do in order to find $\hat{\theta}$, and how many local optima are on the route to $\hat{\theta}$. Consequently, a poor initial starting point may lead to models that take longer to train and/or may lead our learning algorithm to converge to sub-optimal solutions (i.e., an alternative local optimum) (Dodge et al., 2020; Sellam et al., 2022). This can be the case even when only estimating the final layer of a network, e.g., when building a classifier by appending a new layer to a pretrained model—a recent, widely-adopted practice in NLP (Dodge et al., 2020).

Methods for initializing the parameters of neural language models are largely the same as those for initializing other neural networks. Perhaps the simplest approach is to randomly initialize all parameters, e.g., using a uniform or normal random variable generator. The parameters of these generators (mean, standard deviation, bounds) are considered hyperparameters of the learning algorithm. Subsequent methods have iterated on this strategy to develop methods that take into account optimization dynamics or model architectures. One consideration that is particularly relevant for language models is that the input and output sizes of the embedding layer and the fully connected layer can be very different; this exacerbates the problem of vanishing or exploding gradients during training. For example, Glorot, Xavier and Bengio, Yoshua (2010) proposed Xavier init, which keeps the variance of the input and output of all layers within a similar range in order to prevent vanishing or exploding gradients; He et al. (2015) proposed a uniform initialization strategy specifically designed to work with ReLU activation units. Using uniform random variables during parameter initialization can likewise alleviate the problem of vanishing gradients. While most deep learning libraries use thoughtfully-selected initialization strategies for neural networks, it is important to internalize the variance in performance that different strategies can cause.

Early Stopping

As previously discussed, performance on $\mathcal{D}_{\text{train}}$ is not always the best indicator of model performance. Rather, even if our objective continues to increase as we optimize over model parameters, performance on held-out data, i.e., $\mathcal{D}_{\text{test}}$ or even \mathcal{D}_{val} , may suffer as the model starts to overfit to the training data. This phenomenon inspires a practice called **early stopping**, where we stop updating model parameters before reaching

¹³While this logic holds even for samples of size 1 (which is the sample size for standard SGD by definition), basing updates off of single samples can lead to noisy updates. Depending on resource constraints, batch sizes of a few hundred are often used, leading to much more stable training (although in the face of memory constraints, larger batch sizes can be mimicked by accumulating, i.e., averaging, gradients across multiple batches when computing update directions). Batch size itself is often viewed as a model hyperparameter that can have a significant effect on model performance.

(approximately) optimal model parameter values w.r.t. $\mathcal{D}_{\text{train}}$. Instead, we base our stopping criterion C off of model performance on \mathcal{D}_{val} as a quantification of generalization performance, a metric other than that which model parameters are optimized for, or just a general slow down in model improvement on the training objective.

Early stopping sacrifices better training performance for better generalization performance; in this sense, it can also be viewed as a regularization technique, a topic which we discuss next. As with many regularization techniques, early stopping can have adverse effects as well. Recent work suggests that many models may have another period of learning after an initial period of plateauing train/validation set performance. Indeed, a sub-field has recently emerged studying the “grokking” phenomenon (Power et al., 2022), when validation set performance suddenly improves from mediocre to near perfect after a long period in which it appears that model learning has ceased, or even that the model has overfit to the training data. Thus, it is unclear whether early stopping is always a good practice.

3.2.4 Regularization Techniques

Our goal during learning is to produce a model p_{θ} that generalizes beyond the observed data; a model that perfectly fits the training data but produces unrealistic estimates for a new datapoint is of little use. Exactly fitting the empirical distribution is therefore perhaps not an ideal goal. It can lead to **overfitting**, which we informally define as the situation when a model uses spurious relationships between inputs and target variables observed in training data in order to make predictions. While this behavior decreases training loss, it generally harms the model’s ability to perform on unseen data, for which such spurious relationships likely do not hold.

To prevent overfitting, we can apply some form of regularization.

Principle 3.2.3 (Regularization). *Regularization is a modification to a learning algorithm that is intended to increase a model’s generalization performance, perhaps at the cost of training performance.*¹⁴

There are many ways of implementing regularization, such as smoothing a distribution towards a chosen baseline or adding a penalty to the loss function to reflect a prior belief that we may have about the values model parameters should take on Hastie et al. (2001); Bishop (2006). Further, many regularization techniques are formulated for specific model architecture: for example, the count-based smoothing methods used n -gram language models (Ney et al., 1994; Gale and Sampson, 1995). Here we specifically consider the forms of regularization often used in the estimation of neural language models. Most fall into two categories: methods that try to ensure a model’s robustness to yet unseen (or rarely seen) inputs—e.g., by introducing noise into the optimization process—or methods that add a term to our loss function that reflects biases we would like to impart on our model. This is by no means a comprehensive discussion of regularization techniques, for which we refer the reader to Ch.7 of Goodfellow et al. (2016).

Weight Decay

A bias that we may wish to impart on our model is that not all the variables available to the model may be necessary for an accurate prediction. Rather, we hope for our model to learn the simplest mapping from inputs to target variables, as this is likely the function that will be most robust to statistical noise.¹⁵ This bias can be operationalized using regularization techniques such as weight decay (Goodfellow et al., 2016)—also often referred to as ℓ_2 regularization. In short, a penalty for the ℓ_2 norm of θ is added to ℓ . This should in theory discourage the learning algorithm from assigning high values to model parameters corresponding to variables with only a noisy relationship with the output, instead assigning them a value close to 0 that reflects the a non-robust relationship.

¹⁴Adapted from Goodfellow et al. (2016), Ch. 7.

¹⁵This philosophy can be derived from Occam’s Razor, i.e., the principle that one should search for explanations constructed using the smallest possible set of elements.

Entropy Regularization

One sign of overfitting in a language model p_θ is that it places effectively all of its probability mass on a single symbol.¹⁶ Rather, we may want the distributions output by our model to generally have higher entropy, i.e., following the principle of maximum entropy: “the probability distribution which best represents the current state of knowledge about a system is the one with largest entropy, in the context of precisely stated prior data” Jaynes (1957). Several regularization techniques, which we refer to as **entropy regularizers**, explicitly penalize the model for low entropy distributions.

Label smoothing (Szegedy et al., 2015) and the confidence penalty (Pereyra et al., 2017) add terms to ℓ to penalize the model for outputting peaky distributions. Explicitly, label smoothing reassigns a portion of probability mass in the reference distribution from the ground truth symbol to all other symbols in the vocabulary. It is equivalent to adding a term $D_{\text{KL}}(u \parallel p_\theta)$ to ℓ , where u is the uniform distribution. The confidence penalty regularizes against low entropy distribution by adding a term $H(p_\theta)$ to ℓ that encourage high entropy in model outputs. The general class of entropy regularizers have proven effective in training neural models Meister et al. (2020).

Dropout

Regularization also encompasses methods that expose a model to noise that can occur in the data at inference time. The motivation behind such methods is both to penalize a model for being overly dependent on any given variable (whether directly from the input or somewhere further along in the computational graph) for making predictions. Dropout does this explicitly by randomly “dropping” variables from a computation in the network (Srivastava et al., 2014).

More formally, consider a model defined as a series of computational nodes, where any given node is the product of the transformation of previous nodes. When dropout is applied to the module that contains that node, then the node is zeroed out with some percentage chance, i.e., it is excluded in all functions that may make use of it to compute the value of future nodes. In this case, the model will be penalized if it relied completely on the value of that node for any given computation in the network. Dropout can be applied to most variables within a model, e.g., the inputs to the model itself, the inputs to the final linear projection in a feed-forward layer, or the summands in the attention head of a Transformer. Note that at inference time, all nodes are used to compute the model’s output.¹⁷

Batch and Layer Normalization

Rescaling variables within a network helps with training stability, and further, with generalization by keeping variables within the same range and with unit variance. Specifically, batch normalization helps regularize the problem of covariate shift, where the distribution of features (both the input features and the variables corresponding to transformed features within a network) differs between the training data and the data at inference time. Batch normalization alleviates this problem by recentering (around 0) and scaling (such that data points have unit variance) data points such that the data flowing between intermediate layers of the network follows approximately the same distribution between batches. Layer normalization likewise performs centering and rescaling, albeit across features rather than across data points. Specifically, normalization is performed so that all of the feature values within a data point have mean 0 and unit variance.

¹⁶The softmax transformation serves as somewhat of a regularizer against this behavior since it does not allow any symbol be assigned a probability of 0.

¹⁷Some form of renormalization is typically performed to account for the fact that model parameters are learned with only partial availability of variables. Thus when all variables are used in model computations, the scale of the output will (in expectation) be larger than during training, potentially leading to poor estimates.

Chapter 4

Classical Language Models

Next we turn to two classical language modeling frameworks: finite-state language models (a natural generalization of the well-known n -gram models) in §4.1 and pushdown language models §4.2. Although the most successful approaches to language modeling are based on neural networks, the study of older approaches to language modeling is invaluable. First, due to the simplicity of the models, learning how they work helps distill concepts. And, moreover, they often serve as important baselines in modern NLP and provide very useful insights into the capabilities of modern architectures as we will see when we discuss modern architectures in Chapter 5.

In the spirit of our question-motivated investigation, we will focus on the following two questions.

Question 4.0.1. *How can we tractably represent all conditional distributions of the form $p_{SM}(y \mid \mathbf{y})$ in a simple way?*

Question 4.0.2. *How can we tractably represent the hierarchical structure of human language?*

4.1 Finite-state Language Models

After rigorously defining what language models are (and what they are not) and discussing how we can estimate them, it is time to finally introduce our first class of language models—those based on finite-state automata. Language models derived from probabilistic finite-state automata are some of the simplest classes of language models because they definitionally distinguish a *finite* numbers of contexts when modeling the conditional distribution of the next possible symbol $p_M(y \mid \mathbf{y})$. We first give an intuitive definition of a finite-state language model and then introduce a more formal definition, which we will use throughout the rest of the section.

Definition 4.1.1 (Informal definition of a finite-state language model). *A language model p_{LM} is **finite-state** if it defines only finitely many unique conditional distributions $p_{LM}(y \mid \mathbf{y})$. In other words, there are only finitely many contexts \mathbf{y} which define the distribution over the next symbol, $p_{LM}(y \mid \mathbf{y})$.*

Intuitively, this framework might be useful because it *bounds* the number of unique conditional distributions we have to learn. However, as we will see later in this chapter, finite-state language models are not sufficient for modeling human language. Nevertheless, they can still offer a baseline for modeling more complex phenomena. They also offer a useful theoretical tool in the understanding of neural language models, which we will discuss in Chapter 5.

4.1.1 Weighted Finite-state Automata

Before we introduce finite-state *language models*, we go on a brief detour into the theory of finite-state *automata*. As we will see, finite-state automata are a tidy and well-understood formalism. As we will see later in §5.1.7, they also provide a solid and convenient theoretical framework for understanding modern neural language models, e.g., those based on recurrent neural networks and transformers. We, therefore, begin by briefly introducing the theory of finite-state automata with real-valued weights.

Finite-state Automata

In words, finite-state automata are one of the simplest devices for defining a formal language (cf. Definition 2.3.5). We give a formal definition below.

Definition 4.1.2 (Finite-state Automata). A *finite-state automaton* (FSA) is a 5-tuple $(\Sigma, Q, I, F, \delta)$ where

- Σ is an alphabet;
- Q is a finite set of states;
- $I \subseteq Q$ is the set of initial states;
- $F \subseteq Q$ the set of final or accepting states;
- A finite multiset $\delta \subseteq Q \times (\Sigma \cup \{\varepsilon\}) \times Q$.¹ Elements of δ are generally called **transitions**.

The name, *finite-state* automaton, stems from the requirement that the set of states Q is finite, which stands in contrast to the remaining formalisms we will cover in this course, e.g., pushdown automata and recurrent neural networks. We will denote a general finite-state automaton with a (subscripted) \mathcal{A} . We will also adopt a more suggestive notation for transitions by denoting a transition (q_1, a, q_2) as $q_1 \xrightarrow{a} q_2$.

An FSA can be graphically represented as a labeled, directed multi-graph.² The vertices in the graph represent the states $q \in Q$ and the (labeled) edges between them the transitions in δ . The labels on the edges correspond to the input symbols $a \in \Sigma$ which are consumed when transitioning over the edges. The initial states $q_i \in I$ are marked by a special incoming arrow while the final states $q_f \in F$ are indicated using a *double* circle.

Example 4.1.1. An example of an FSA can be seen in Fig. 4.1. Formally, we can specify it as

- $\Sigma = \{a, b, c\}$
- $Q = \{1, 2, 3\}$
- $I = \{1\}$
- $F = \{3\}$
- $\delta = \{(1, a, 2), (1, b, 3), (2, b, 2), (2, c, 3)\}$

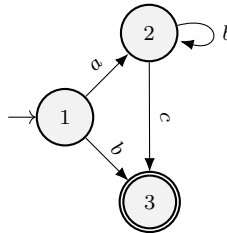


Figure 4.1: Example of a simple FSA.

A finite-state automaton sequentially reads in individual symbols of an **input string** $y \in \Sigma^*$ and transitions from state to state according to the transition function δ . The traversal through the automaton starts in a state $q_i \in I$ (more precisely, it acts as if starting from all of them in parallel). It then transitions from state q into the state q' upon reading the symbol a if and only if $q \xrightarrow{a} q' \in \delta$. ε -labeled transitions,

¹The fact that it is a multiset reflects that it can contain multiple copies of the same element (i.e., transitions between the same pair of states with the same symbol).

²The *multi*- aspect of the multi-graph refers to the fact that we can have multiple transitions from any pair of states and labeled refers to the fact that we label those transitions with symbols from the alphabet Σ .

however, allow a finite-state machine to transition to a new state without consuming a symbol. This is in line with ε 's definition as an empty string.

A natural question to ask at this point is what happens if for a state–symbol pair (q, a) there is *more than one* possible transition allowed under the relation δ . In such a case, we take *all* implicit transitions simultaneously, which leads us to a pair of definitions.

Definition 4.1.3. A FSA $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ is **deterministic** if

- it does not have any ε -transitions;
- for every $(q, a) \in Q \times \Sigma$, there is at most one $q' \in Q$ such that $q \xrightarrow{a} q' \in \delta$;
- there is a single initial state, i.e., $|I| = 1$.

Otherwise, \mathcal{A} is **non-deterministic**.

An important, and perhaps not entirely obvious, result is that the classes of deterministic and non-deterministic FSA are *equivalent*, in the sense that you can always represent a member of one class with a member of the other.

If the automaton ends up, after reading in the last symbol of the input string, in one of the final states $q_\varphi \in F$, we say that the automaton **accepts** that string. A finite-state automaton is therefore a computational device that determines whether a string satisfies a condition (namely, the condition that the automaton, by starting in an initial state and following one of the paths labeled with that string, ends in a final state). A string that satisfies this condition is said to be *recognized* by the automaton and the set of all strings satisfying this condition form the *language* of the automaton.³

Definition 4.1.4. Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be an finite-state automaton. The **language** of \mathcal{A} , $L(\mathcal{A})$ is defined as

$$L(\mathcal{A}) \stackrel{\text{def}}{=} \{y \mid y \text{ is recognized by } \mathcal{A}\} \quad (4.1)$$

Abstractly, a finite-state automaton is hence a specification of a set of *rules* that strings must satisfy to be included in its language. The set of languages that finite-state automata can recognize is known as the class of **regular languages**.

Definition 4.1.5. A language $L \subseteq \Sigma^*$ is **regular** if and only if it can be recognized by an unweighted finite-state automaton, i.e., if there exists a finite-state automaton \mathcal{A} such that $L = L(\mathcal{A})$.

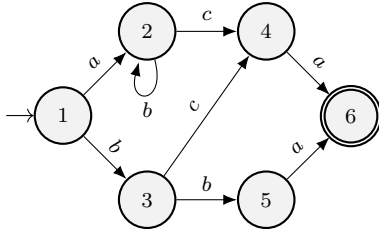
Example 4.1.2. Additional simple examples of FSAs are shown in Fig. 4.2. The FSA in Fig. 4.2a, for example, can formally be defined with

- $\Sigma = \{a, b, c\}$
- $Q = \{1, 2, 3, 4, 5, 6\}$
- $I = \{1\}$
- $F = \{6\}$
- $\delta = \{(1, a, 2), (1, b, 3), (2, b, 2), (2, c, 4), (3, c, 4), (3, b, 5), (4, a, 6), (5, a, 6)\}$

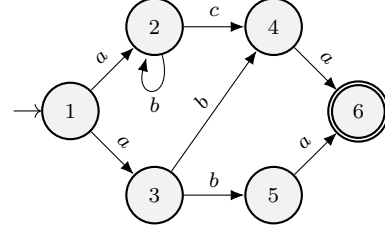
The FSA in Fig. 4.2a is deterministic while the one in Fig. 4.2b is non-deterministic. A few examples of strings accepted by the \mathcal{A}_1 include *bba*, *bca*, *aca*, *abca*, *abbca*, *abbbca*, \dots . In fact, due to the self-loop at state 2, the symbol *b* can appear an arbitrary number of times at position 2 in the accepted string *abca*. Notice that, starting from the state 1 and following the transitions dictated by any of the accepted strings, we always end up in the only final state, state 6. In particular, the string “*abbca*” is accepted with the following set of transitions in \mathcal{A}_1 :

$$1 \xrightarrow{a} 2, 2 \xrightarrow{b} 2, 2 \xrightarrow{b} 2, 2 \xrightarrow{c} 4, 4 \xrightarrow{a} 6.$$

³We also say that the automaton *recognizes* this set of strings (language).



(a) A deterministic FSA, \mathcal{A}_1 . Each state only has one outgoing transition labeled with the same symbol.



(b) A non-deterministic FSA, \mathcal{A}_2 . State 1 has two outgoing transitions labeled with a whereas state 3 has two outgoing transitions labeled with b .

Figure 4.2: Examples of a deterministic and a non-deterministic FSA.

Weighted Finite-state Automata

A common and very useful augmentation to finite-state automata is through the addition of *weights* on the transitions. The general theory of weighted automata makes use of semiring theory, which is beyond the scope of this course.⁴ In this course, we will limit ourselves to the study of automata with real-valued weights.

Definition 4.1.6 (Real-weighted Finite-State Automaton). *A **real-weighted finite-state automaton** (WFSA) \mathcal{A} is a 5-tuple $(\Sigma, Q, \delta, \lambda, \rho)$ where*

- Σ is a finite alphabet;
- Q is a finite set of states;
- $\delta \subseteq Q \times (\Sigma \cup \{\varepsilon\}) \times \mathbb{R} \times Q$ a finite multiset of transitions;⁵
- $\lambda : Q \rightarrow \mathbb{R}$ a weighting function over Q ;
- $\rho : Q \rightarrow \mathbb{R}$ a weighting function over Q .

Notice that we omit the initial and final state sets from the definition of WFSA. Those can implicitly be specified by the states given non-zero initial or final weights by the λ and ρ functions, i.e., $I = \{q \in Q \mid \lambda(q) \neq 0\}$ and $F = \{q \in Q \mid \rho(q) \neq 0\}$. We might refer to them in the text later for notational convenience and clarity of exposition. We will also sometimes denote transition weights with $\omega \left(q \xrightarrow{a/w} q' \right) \stackrel{\text{def}}{=} w$.

Graphically, we write the transition weights on the edges of the graph representing the WFSA after the output symbol, separated by a “/”. The same separator is also used to separate the state name from its initial/final weights, which are written in the node.

Example 4.1.3. Fig. 4.3 shows a weighted version of the FSA from Fig. 4.2a above.

The connection of WFSA to graphs makes it natural to define a set of transition matrices specified by a WFSA.

Definition 4.1.7. Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a WFSA. For any $a \in \Sigma$, we define the **symbol-specific transition matrix** $\mathbf{T}^{(a)}$ as the transition matrix of the graph restricted to a -labeled transitions. We also define the (full) **transition matrix** as $\mathbf{T} \stackrel{\text{def}}{=} \sum_{a \in \Sigma} \mathbf{T}^{(a)}$.

⁴Semirings and semiring-weighted formal languages are covered in detail in the Advanced Formal Language Theory course offered at ETH as well.

⁵Again, we use the notation $q \xrightarrow{a/w} q'$ to denote $(q, a, w, q') \in \delta$.

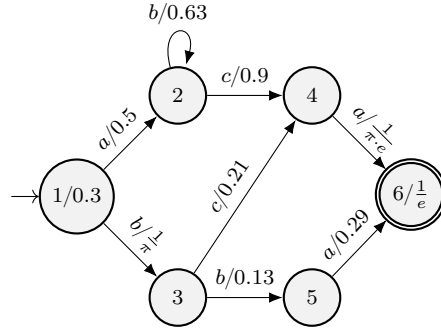


Figure 4.3: The WFSA corresponding to the FSA from Fig. 4.2a.

Example 4.1.4. Consider the WFSA \mathcal{A} in Fig. 4.3. The (symbol-specific) transition matrices for \mathcal{A} are

$$\mathbf{T}^{(a)} = \begin{pmatrix} 0 & 0.5 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \frac{1}{\pi \cdot e} \\ 0 & 0 & 0 & 0 & 0 & 0.29 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\mathbf{T}^{(b)} = \begin{pmatrix} 0 & 0 & \frac{1}{\pi} & 0 & 0 & 0 \\ 0 & 0.63 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.13 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\mathbf{T}^{(c)} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.9 & 0 & 0 \\ 0 & 0 & 0 & 0.21 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\mathbf{T} = \mathbf{T}^{(a)} + \mathbf{T}^{(b)} + \mathbf{T}^{(c)} = \begin{pmatrix} 0 & 0.5 & \frac{1}{\pi} & 0 & 0 & 0 \\ 0 & 0.63 & 0 & 0.9 & 0 & 0 \\ 0 & 0 & 0 & 0.21 & 0.13 & 0 \\ 0 & 0 & 0 & 0 & 0 & \frac{1}{\pi \cdot e} \\ 0 & 0 & 0 & 0 & 0 & 0.29 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

Paths and Path Weights

A path is an important concept when talking about (weighted) finite-state automata as it defines the basic structure by which a string is recognized or weighted. We now give a formal definition of a path and discuss how to weight paths.

Definition 4.1.8 (Path). A **path** π is an element of δ^* with consecutive transitions, meaning that it is of the form $\left(q_1 \xrightarrow{\bullet/\bullet} q_2, q_2 \xrightarrow{\bullet/\bullet} q_3 \cdots q_{n-1} \xrightarrow{\bullet/\bullet} q_n \right)$, where \bullet is a placeholder.⁶ The **length** of a path is the number of transition in it; we denote the length as $|\pi|$. We use $p(\pi)$ and $n(\pi)$ to denote the origin and the destination of a path, respectively. The **yield** of a path is the concatenation of the input symbols on the

⁶Notice we use the Kleene closure on the set δ here. It thus represents any sequence of transitions $\in \delta$

edges along the path, which we will mark with $\mathbf{s}(\pi)$. Furthermore, we denote sets of paths with capital Π . Throughout the text, we will use a few different variants involving Π to avoid clutter:

- $\Pi(\mathcal{A})$ as the set of all paths in automaton \mathcal{A} ;
- $\Pi(\mathcal{A}, \mathbf{y})$ as the set of all paths in automaton \mathcal{A} with yield $\mathbf{y} \in \Sigma^*$;
- $\Pi(\mathcal{A}, q, q')$ as the set of all paths in automaton \mathcal{A} from state q to state q' .

One of the most important questions when talking about weighted formalisms like weighted finite-state automata is how to *combine* weights of atomic units like transitions into weights of complete *structures*.⁷ We begin by multiplicatively combining the weights of individual transitions in a path into the weights of the full path.

Definition 4.1.9 (Path Weight). The **inner path weight** $w_I(\pi)$ of a path $\pi = q_1 \xrightarrow{a_1/w_1} q_2 \cdots q_{N-1} \xrightarrow{a_N/w_N} q_N$ is defined as

$$w_I(\pi) = \prod_{n=1}^N w_n. \quad (4.2)$$

The **(full) path weight** of the path π is then defined as

$$w(\pi) = \lambda(p(\pi)) w_I(\pi) \rho(n(\pi)). \quad (4.3)$$

A path π is called **accepting** or **successful** if $w(\pi) \neq 0$.

The inner path weight is therefore the product of the weights of the transitions on the path, while the (full) path weight is the product of the transition weights as well as the initial and final weights of the origin and the destination of the path, respectively.

String Acceptance Weights and Weighted Regular Languages

When we introduced unweighted finite-state automata, we defined the important concept of recognizing a string and recognizing a language. We generalize these concepts to the very natural quantity of the weight assigned by a WFSa to a string $\mathbf{y} \in \Sigma^*$, i.e., its acceptance weight, or stringsum, as the sum of the weights of the paths that yield \mathbf{y} .

Definition 4.1.10 (Strings sum). The **stringsum**, **string weight**, or **acceptance weight** of a string $\mathbf{y} \in \Sigma^*$ under a WFSa \mathcal{A} is defined as

$$\mathcal{A}(\mathbf{y}) \stackrel{\text{def}}{=} \sum_{\pi \in \Pi(\mathcal{A}, \mathbf{y})} w(\pi). \quad (4.4)$$

This naturally generalizes the notion of acceptance by an unweighted FSA—whereas an unweighted FSA only makes a binary decision of accepting or rejecting a string, a weighted FSA always accepts a string with a specific weight. This leads to the definition of the *weighted language* of the WFSa.

Definition 4.1.11. Let \mathcal{A} be a WFSa. Its **(weighted) language** is defined as

$$L(\mathcal{A}) \stackrel{\text{def}}{=} \{(\mathbf{y}, \mathcal{A}(\mathbf{y})) \mid \mathbf{y} \in \Sigma^*\}. \quad (4.5)$$

We say a language is a weighted regular language if it is a language of some WFSa:

Definition 4.1.12. A **weighted language** L is a **weighted regular language** if there exists a WFSa \mathcal{A} such that $L = L(\mathcal{A})$.

Lastly, we also define the full and state-specific allsum of the automaton. The former refers to the total weight assigned to *all* possible strings, or all possible paths whereas the latter refers to the sum of the path weights of the paths stemming from a specific state.

⁷In the case of WFSAs, a structure is a path. In the next section, we will see how to combine weights from basic units into trees.

Definition 4.1.13 (State-specific allsum). *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a WFSA. The **allsum** of a state $q \in Q$ is defined as*

$$Z(\mathcal{A}, q) = \sum_{\substack{\pi \in \Pi(\mathcal{A}) \\ q_1 = q}} w_I(\pi) \rho(n(\pi)). \quad (4.6)$$

State-specific allsums are also referred to as the **backward values** in the literature and are often denoted as $\beta(q)$.

Definition 4.1.14 (WFSA allsum). *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a WFSA. The **allsum** of \mathcal{A} is defined as*

$$Z(\mathcal{A}) = \sum_{y \in \Sigma^*} \mathcal{A}(y) = \sum_{y \in \Sigma^*} \sum_{\pi \in \Pi(\mathcal{A}, y)} w(\pi) = \sum_{\pi \in \Pi(\mathcal{A})} w(\pi). \quad (4.7)$$

The second equality in Eq. (4.7) comes from the crucial observation that the double sum in the second term sums over precisely *all paths* of the automaton \mathcal{A} , which is where the name of the quantity comes from *allsum*.⁸ This is easy to see if we consider that by summing over all possible strings, we enumerate all possible path yields, and each path in the automaton has a yield $\in \Sigma^*$. $Z(\mathcal{A})$ is again the result of summing over infinitely many terms (whether the set of strings in Σ^* of the infinitely many paths in a cyclic WFSA), and might therefore not necessarily be finite. For reasons which will become clear shortly, we will say that a WFSA \mathcal{A} is **normalizable** if $Z(\mathcal{A}) < \infty$.

Note that the sum in Eq. (4.4) only contains one term if the automaton is deterministic. Whenever the automaton is non-deterministic, or when we are interested in the sum of paths with *different* yields as in Eq. (4.7), the interactions (namely, the *distributive law*) between the *sum* over the different *paths* and the *multiplications* over the *transitions* in the paths play an important role when designing efficient algorithms. Indeed, many algorithms defined for WFSAs rely on decompositions of such sums enabled by the distributive law.⁹

Accessibility and Probabilistic Weighted Finite-state Automata

An important property of states of a WFSA which we will need when investigating the tightness of finite-state language models is accessibility.

Definition 4.1.15. *A state $q \in Q$ of a WFSA is **accessible** if there is a non-zero-weighted path to q from some state q_i with $\lambda(q_i) \neq 0$; it is **co-accessible state** if there is a non-zero-weighted path from q to some state q_f with $\rho(q_f) \neq 0$. It is **useful** if it is both accessible and co-accessible, i.e., q appears on some non-zero-weighted accepting path.*

Definition 4.1.16. *Trimming a WFSA means removing its useless states.¹⁰ Removing the non-useful states means removing their rows and columns from \mathbf{T} as well as their rows from $\vec{\lambda}$ and $\vec{\rho}$, yielding possibly smaller \mathbf{T}' , $\vec{\lambda}'$ and $\vec{\rho}'$.*

We will use WFSAs to specify language models. However, not every WFSA is a language model, i.e., a distribution over strings. Generally, the weight of a string could be negative if we allow arbitrary real weights. Thus, a restriction we will impose on *all* weighted automata that represent finite-state language models is that the weights be *non-negative*.

Furthermore, a special class of WFSAs that will be of particular interest later is probabilistic WFSAs.

Definition 4.1.17 (Probabilistic Weighted Finite-State Automaton). *A WFSA $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ is **probabilistic** (a PFSA) if*

$$\sum_{q \in Q} \lambda(q) = 1 \quad (4.8)$$

⁸Analogously, given some (implicitly defined) set of paths \mathcal{S} , we will name the sum over the weights of the paths in \mathcal{S} the allsum over \mathcal{S}

⁹Many such examples are covered in the Advanced Formal Language Theory course.

¹⁰This does not affect the weights of the strings with $w(y) \neq 0$.

and, for all $q \in Q$ and all outgoing transitions $q \xrightarrow{a/w} q' \in \delta$ it holds that

$$\lambda(q) \geq 0 \quad (4.9)$$

$$\rho(q) \geq 0 \quad (4.10)$$

$$w \geq 0 \quad (4.11)$$

and

$$\sum_{q \xrightarrow{a/w} q'} w + \rho(q) = 1. \quad (4.12)$$

This means that the initial weights of all the states of the automaton form a probability distribution (the initial weight of a state corresponds to the *probability* of starting in it), as well as that, for any state q in the WSFA, the weights of its outgoing transitions (with any label) together with its final weight form a valid discrete probability distribution. In a certain way, probabilistic finite-state automata naturally correspond to *locally normalized* language models, as we explore in the next subsection.

The eos symbol and the final weights. Notice that the final weights in a PFSA play an analogous role to the EOS symbol: the probability of ending a path in a specific state q —and therefore ending a string—is q ’s final weight! That is, the probability $\rho(q_\varphi)$ for some $q_\varphi \in Q$, representing the probability of ending the path in q_φ , is analogous to the probability of ending a string \mathbf{y} , $p_{\text{SM}}(\text{EOS} \mid \mathbf{y})$, where q_φ “represents” the string (history) \mathbf{y} .¹¹ When modeling language with weighted finite-state automata, we will therefore be able to avoid the need to specify the special symbol and rather rely on the final weights, which are naturally part of the framework.

4.1.2 Finite-state Language Models

We can now formally define what it means for a language model to be finite-state:

Definition 4.1.18 (Finite-state language models). *A language model p_{LM} is **finite-state** if it can be represented by a weighted finite-state automaton, i.e., if there exists a WFSA $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ such that $L(\mathcal{A}) = L(p_{LM})$. Equivalently, we could say that p_{LM} is finite-state if its language is a weighted regular language.*

On the other hand, given a WFSA \mathcal{A} , there are two established ways of defining a *probability* of string.

String Probabilities in a Probabilistic Finite-state Automaton

In a probabilistic FSA (cf. Definition 4.1.17), any action from a state $q \in Q$ is associated with a probability. Since the current state completely encodes all the information of the input seen so far in a finite-state automaton, it is intuitive to see those probabilities as conditional probabilities of the next symbol given the input seen so far. One can, therefore, define the probability of a path as the product of these individual “conditional” probabilities.

Definition 4.1.19 (Path probability in a PFSA). *We call the weight of a path $\pi \in \Pi(\mathcal{A})$ in a probabilistic FSA the **probability** of the path π .*

This alone is not enough to define the probability of any particular string $\mathbf{y} \in \Sigma^*$ since there might be multiple accepting paths for \mathbf{y} . Naturally, we define the probability of \mathbf{y} as the sum of the individual paths that recognize it:

Definition 4.1.20 (String probability in a PFSA). *We call the stringsum of a string $\mathbf{y} \in \Sigma^*$ in a probabilistic FSA the **probability** of the string \mathbf{y} :*

$$p_{\mathcal{A}}(\mathbf{y}) \stackrel{\text{def}}{=} \mathcal{A}(\mathbf{y}). \quad (4.13)$$

¹¹Due to the possible non-determinism of WFSAs, the connection is of course not completely straightforward, but the point still stands.

Crucially, notice that these two definitions did not require any *normalization* over all possible paths or strings. This closely resembles the way we defined locally normalized models based on the conditional probabilities of a sequence model. Again, such definitions of string probabilities are attractive as the summation over all possible strings is avoided. However, a careful reader might then ask themselves: do these probabilities actually sum to 1, i.e., is a probabilistic FSA *tight*? As you might guess, they might *not*.¹² We explore this question in §4.1.4.

String Probabilities in a General Weighted Finite-state Automaton

To define string probabilities in a general weighted FSA, we use the introduced notions of the stringsum and the allsum. The allsum allows us to tractably *normalize* the stringsum to define the globally normalized probability of a string \mathbf{y} as the *proportion* of the total weight assigned to all strings that is assigned to \mathbf{y} .¹³

Definition 4.1.21 (String probability in a WFSA). *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a normalizable WFSA with non-negative weights. We define the probability of a string $\mathbf{y} \in \Sigma^*$ under \mathcal{A} as*

$$p_{\mathcal{A}}(\mathbf{y}) \stackrel{\text{def}}{=} \frac{\mathcal{A}(\mathbf{y})}{Z(\mathcal{A})}. \quad (4.14)$$

Language Models Induced by a WFSA

With the notions of string probabilities in both probabilistic and general weighted FSAs, we can now define the language model induced by \mathcal{A} as follows.

Definition 4.1.22 (A language model induced by a WFSA). *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a WFSA. We define the **language model induced by \mathcal{A}** as the following probability distribution over Σ^**

$$p_{LMA}(\mathbf{y}) \stackrel{\text{def}}{=} p_{\mathcal{A}}(\mathbf{y}). \quad (4.15)$$

It is easy to see that while global normalization requires the computation of the allsum, language models induced by weighted FSAs through Eq. (4.14) are *globally normalized* and thus always tight. In the next subsection, we consider *how* the quantities needed for computing Eq. (4.14) can be computed. Of particular interest will be the quantity $Z(\mathcal{A})$, as it involves the summation over possibly infinitely many terms and therefore requires some clever tricks to be computed.

4.1.3 Normalizing Finite-state Language Models

In this subsection, we develop an algorithm for normalizing a globally normalized language model (cf. Definition 2.4.2) defined by a WFSA, i.e., an algorithm for computing the allsum $Z(\mathcal{A})$ whenever this quantity is finite. Moreover, the derivation will also reveal necessary and sufficient conditions for FSAs to be normalizable.

Converting a matrix of pairwise pathsums to the allsum. Before we consider how to compute $Z(\mathcal{A})$, let us first consider a much simpler problem. Suppose we had a matrix \mathbf{M} , which contained at the entry \mathbf{M}_{ij} the sum of all the inner weights over all paths between the states i and j , i.e.,

$$\mathbf{M}_{ij} = \sum_{\pi \in \Pi(\mathcal{A}, i, j)} w_{\mathbf{I}}(\pi).$$

¹²Notice that, however, whenever a PFSA is tight, its allsum is 1.

¹³We will see how the allsum can be computed tractably in §4.1.3.

How could we then compute the quantity $Z(\mathcal{A})$?

$$Z(\mathcal{A}) = \sum_{\pi \in \Pi(\mathcal{A})} w(\pi) \quad (4.16)$$

$$= \sum_{\pi \in \Pi(\mathcal{A})} \lambda(p(\pi)) w_I(\pi) \rho(n(\pi)) \quad (4.17)$$

$$= \sum_{i,j \in Q} \sum_{\pi \in \Pi(\mathcal{A}, i, j)} \lambda(p(\pi)) w_I(\pi) \rho(n(\pi)) \quad (4.18)$$

$$= \sum_{i,j \in Q} \sum_{\pi \in \Pi(\mathcal{A}, i, j)} \lambda(i) w_I(\pi) \rho(j) \quad (4.19)$$

$$= \sum_{i,j \in Q} \lambda(i) \left(\sum_{\pi \in \Pi(\mathcal{A}, i, j)} w_I(\pi) \right) \rho(j) \quad (4.20)$$

$$= \sum_{i,j \in Q} \lambda(i) \mathbf{M}_{ij} \rho(j) \quad (4.21)$$

$$= \vec{\lambda} \mathbf{M} \vec{\rho}, \quad (4.22)$$

where $\vec{\lambda}$ and $\vec{\rho}$ denote the vectors resulting from the “vectorization” of the functions λ and ρ , i.e., $\vec{\lambda}_n = \lambda(n)$ and $\vec{\rho}_n = \rho(n)$. This also explains the naming of the functions λ and ρ : the initial weights function λ , “lambda” appears on the left side of the closed form expression for $Z(\mathcal{A})$ and the definition of the path weight (cf. Eq. (4.3)), whereas the final weights function ρ , rho, appears on the right side of the expression and the definition of the path weight.

Computing the matrix of pairwise pathsums. Let \mathbf{T} be the transition matrix of the automaton \mathcal{A} . Notice that the entry \mathbf{T}_{ij} by definition contains the sum of the inner weights of all paths of length exactly 1 (individual transitions) between the states i and j . We also define $\mathbf{T}^0 = \mathbf{I}$, meaning that the sum of the weights of the paths between i and j of length zero is 0 if $i \neq j$ and 1 (the unit for multiplication) if $i = j$. This corresponds to not transitioning, i.e., staying in place, if $i = j$. We next state a basic result from graph theory.

Lemma 4.1.1. *Let \mathbf{T} be the transition matrix of some weighted directed graph \mathcal{G} . Then the matrix \mathbf{T}^d contains the allsum of all paths of length exactly d , i.e.,*

$$\mathbf{T}_{i,j}^d = \sum_{\substack{\pi \in \Pi(\mathcal{A}, i, j) \\ |\pi|=d}} w_I(\pi). \quad (4.23)$$

Proof. By induction on the path length. Left as an exercise for the reader. ■

It follows directly that the matrix

$$\mathbf{T}^{\leq d} \stackrel{\text{def}}{=} \sum_{k=1}^d \mathbf{T}^k$$

contains the pairwise pathsums of paths of length *at most* d .

In general, the WFSA representing a n -gram language model can of course be cyclic. This means that the number of paths in $\Pi(\mathcal{A})$ might be infinite and they might be of arbitrary length (which is the result of looping in a cycle arbitrarily many times). To compute the pairwise pathsums over *all* possible paths, we, therefore, have to compute

$$\mathbf{T}^* \stackrel{\text{def}}{=} \lim_{d \rightarrow \infty} \mathbf{T}^{\leq d} = \sum_{d=0}^{\infty} \mathbf{T}^d. \quad (4.24)$$

This is exactly the matrix form of the *geometric sum*. Similarly to the scalar version, we can manipulate the expression Eq. (4.24) to arrive to a closed-form expression for computing it:

$$\mathbf{T}^* = \sum_{d=0}^{\infty} \mathbf{T}^d \quad (4.25)$$

$$= \mathbf{I} + \sum_{d=1}^{\infty} \mathbf{T}^d \quad (4.26)$$

$$= \mathbf{I} + \sum_{d=1}^{\infty} \mathbf{T} \mathbf{T}^{d-1} \quad (4.27)$$

$$= \mathbf{I} + \mathbf{T} \sum_{d=1}^{\infty} \mathbf{T}^{d-1} \quad (4.28)$$

$$= \mathbf{I} + \mathbf{T} \sum_{d=0}^{\infty} \mathbf{T}^d \quad (4.29)$$

$$= \mathbf{I} + \mathbf{T} \mathbf{T}^*. \quad (4.30)$$

If the inverse of $(\mathbf{I} - \mathbf{T})$ exists, we can further rearrange this equation to arrive at

$$\mathbf{T}^* = \mathbf{I} + \mathbf{T} \mathbf{T}^* \quad (4.31)$$

$$\mathbf{T}^* - \mathbf{T} \mathbf{T}^* = \mathbf{I} \quad (4.32)$$

$$\mathbf{T}^* - \mathbf{T}^* \mathbf{T} = \mathbf{I} \quad (4.33)$$

$$\mathbf{T}^* (\mathbf{I} - \mathbf{T}) = \mathbf{I} \quad (4.34)$$

$$\mathbf{T}^* = (\mathbf{I} - \mathbf{T})^{-1}. \quad (4.35)$$

This means that, if $(\mathbf{I} - \mathbf{T})$ exists, we can compute the pairwise pathsums by simply inverting it! Using the remark above on how to convert a matrix of pairwise pathsums into the full allsum, we can therefore see that we can globally normalize an n -gram language model by computing a matrix inversion! Since the runtime of inverting a $N \times N$ matrix is $\mathcal{O}(N^3)$, and $N = |Q|$ for a transition matrix of a WFSA with states Q , we can globally normalize a n -gram language model in time cubic in the number of its states. This is a special case of the general algorithm by [Lehmann \(1977\)](#). Note, however, that this might still be prohibitively expensive: as we saw, the number of states in a n -gram model grows exponentially with n , and even small n 's and reasonable alphabet sizes might result in a non-tractable number of states in the WFSA with the cubic runtime.

We still have to determine when the infinite sum in Eq. (4.24) converges. One can see by writing out the product \mathbf{T}^d in terms of its eigenvalues that the entries of \mathbf{T}^d diverge towards $\pm\infty$ as soon as the magnitude of any of \mathbf{T} 's eigenvalues is larger than 1. This means that $\|\mathbf{T}\|_2 < 1$ (spectral norm) is a necessary condition for the infinite sum to exist. This is, however, also a *sufficient* condition: if $\|\mathbf{T}\|_2 < 1$, all of \mathbf{T} 's eigenvalues are smaller than 1 in magnitude, meaning that the eigenvalues of $\mathbf{I} - \mathbf{T}$ are strictly positive and the matrix $\mathbf{I} - \mathbf{T}$ is invertible.¹⁴

Speed-ups of the Allsum Algorithm

The introduced algorithm for computing the allsum in a WFSA can, therefore, be implemented as a matrix inverse. This means that its runtime is $\mathcal{O}(|Q|^3)$, which can be relatively expensive. Fortunately, faster algorithms exist for WSAs with more structure (in their transition functions)—for example, the allsum can be computed in time *linear* in the number of transitions if the automaton is acyclic using a variant of the Viterbi algorithm ([Eisner, 2016](#)). Furthermore, if the automaton “decomposes” into many smaller strongly connected components (i.e., subgraphs that are cyclic), but the components are connected sparsely and form an acyclic graph of components, the allsum can also be computed more efficiently using a combination of the

¹⁴ $1 - \lambda$ is an eigenvalue of $\mathbf{I} - \mathbf{T}$ iff λ is an eigenvalue of \mathbf{T} .

algorithms described above and the algorithm for acyclic WFSA, resulting in a possibly large speedup over the original algorithm.

Importantly, the allsum algorithm and all the speed-ups are *differentiable*, meaning that they can be used during the *gradient-based training* (cf. §3.2.3) of a finite-state language model, where the weights are parametrized using some learnable parameters—we will return to this point shortly.

Locally Normalizing a Globally Normalized Finite-state Language Model

As shown in Theorem 2.4.2, any language model (and thus, any globally-normalized model with a normalizable energy function) can also be locally normalized. In the case of finite-state language models, we can actually explicitly construct the WFSA representing the locally normalized variant using a procedure that is conceptually similar to the allsum algorithm described here. In contrast to the procedure we presented here, however, the local normalization algorithm computes the pathsums of the paths stemming from every possible state q *individually* and then “reweights” the *transitions* depending on the pathsums of their target states r . You can think of this as computing the contributions to the entire allsum from q made by all the individual outgoing transitions from q and then normalizing those contributions. This is an instance of the more general **weight pushing** algorithm.¹⁵ This can be summarized by the following theorem:

Theorem 4.1.1 (PFSA and WFSA are equally expressive). *Normalizable weighted finite-state automata with non-negative weights and tight probabilistic finite-state automata are equally expressive.*

In the proof of this theorem, we will make use of the following lemma.

Lemma 4.1.2. *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ and $q \in Q$. Then*

$$Z(\mathcal{A}, q) = \sum_{q \xrightarrow{a/w} q' \in \delta_{\mathcal{A}_L}} \omega\left(q \xrightarrow{a/\cdot} q'\right) Z(\mathcal{A}, q') + \rho(q) \quad (4.36)$$

Proof. You are asked to show this in Exercise 4.1. ■

We can now prove Theorem 4.1.1

Proof. To prove the theorem, we have to show that any WFSA can be written as a PFSA and vice versa.¹⁶

\Leftarrow Since any tight probabilistic FSA is simply a WFSA with $Z(\mathcal{A}) = 1$, this holds trivially.

\Rightarrow Local normalization is a general property of automata resulting from weight pushing. Here, we describe the construction in the special case of working with real-valued weights. See Mohri et al. (2008) for a general treatment.

Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a normalizable WFSA with non-negative weights. We now show that, for any WFSA, there exists a PFSA encoding the same language model. Let $\mathcal{A}_G = (\Sigma, Q, \delta, \lambda, \rho)$ be a trim WFSA that encodes a distribution over Σ^* using Eq. (4.14). We now construct a tight probabilistic finite-state automaton $\mathcal{A}_L = (\Sigma, Q, \delta_{\mathcal{A}_L}, \lambda_{\mathcal{A}_L}, \rho_{\mathcal{A}_L})$ whose language is identical. We define the initial and final weights of the probabilistic FSA as follows.

$$\lambda_{\mathcal{A}_L}(q) \stackrel{\text{def}}{=} \lambda(q) \frac{Z(\mathcal{A}, q)}{Z(\mathcal{A})} \quad (4.37)$$

$$\rho_{\mathcal{A}_L}(q) \stackrel{\text{def}}{=} \frac{\rho(q)}{Z(\mathcal{A}, q)} \quad (4.38)$$

We define the transitions of the probabilistic FSA as follows.

$$\omega_{\mathcal{A}_L}\left(q \xrightarrow{a/\cdot} q'\right) \stackrel{\text{def}}{=} \frac{\omega\left(q \xrightarrow{a/\cdot} q'\right) Z(\mathcal{A}, q')}{Z(\mathcal{A}, q)} \quad (4.39)$$

¹⁵See Mohri et al. (2008) for a more thorough discussion of weight pushing.

¹⁶By “written as”, we mean that the weighted language is the same.

This means that \mathcal{A}_L contains the same transitions as \mathcal{A} , they are simply reweighted. Note that the assumption that \mathcal{A} is trimmed means that all the quantities in the denominators are non-zero.

It is easy to see that the weights defined this way are non-negative due to the non-negativity of \mathcal{A} 's weights. Furthermore, the weights of all outgoing arcs from any $q \in Q$ and its final weight sum to 1:

$$\sum_{q \xrightarrow{a/w} q' \in \delta_{\mathcal{A}_L}} w + \rho_{\mathcal{A}_L}(q) \quad (4.40)$$

$$= \sum_{q \xrightarrow{a/w} q' \in \delta_{\mathcal{A}_L}} \frac{\omega\left(q \xrightarrow{a/\cdot} q'\right) Z(\mathcal{A}, q')}{Z(\mathcal{A}, q)} + \frac{\rho(q)}{Z(\mathcal{A}, q)} \quad (\text{definition of } \delta_{\mathcal{A}_L}) \quad (4.41)$$

$$= \frac{1}{Z(\mathcal{A}, q)} \left(\sum_{q \xrightarrow{a/w} q' \in \delta_{\mathcal{A}_L}} \omega\left(q \xrightarrow{a/\cdot} q'\right) Z(\mathcal{A}, q') + \rho(q) \right) \quad (4.42)$$

$$= 1 \quad (\text{Lemma 4.1.2}) \quad (4.43)$$

It is also easy to see that the initial weights form a probability distribution over the states of the constructed automaton.

$$\sum_{q \in Q} \lambda_{\mathcal{A}_L}(q) = \sum_{q \in Q} \lambda(q) \frac{Z(\mathcal{A}, q)}{Z(\mathcal{A})} \quad (4.44)$$

$$= \frac{1}{Z(\mathcal{A})} \sum_{q \in Q} \lambda(q) Z(\mathcal{A}, q) \quad (4.45)$$

$$= \frac{1}{Z(\mathcal{A})} Z(\mathcal{A}) = 1 \quad (4.46)$$

We now have to show that the probabilities assigned by these two automata match. We will do that by showing that the probabilities assigned to individual *paths* match, implying that stringsums match as well. The probability of a path is defined analogously to a probability of a string, i.e., $p_{\mathcal{A}}(\boldsymbol{\pi}) = \frac{w(\boldsymbol{\pi})}{Z(\mathcal{A})}$ (where $Z(\mathcal{A}) = 1$ for tight probabilistic FSAs). Let then $\boldsymbol{\pi} = \left(q_1 \xrightarrow{a_1/w_1} q_2, \dots, q_{N-1} \xrightarrow{a_{N-1}/w_{N-1}} q_N \right) \in \Pi(\mathcal{A}) = \Pi(\mathcal{A}_L)$. Then, by the definitions of $\omega_{\mathcal{A}_L}$, $\lambda_{\mathcal{A}_L}$, and $\rho_{\mathcal{A}_L}$

$$p_{\mathcal{A}_L}(\boldsymbol{\pi}) = \lambda_{\mathcal{A}_L}(q_1) \left(\prod_{n=1}^{N-1} w_n \right) \rho_{\mathcal{A}_L}(q_N) \quad (4.47)$$

$$= \lambda(q_1) \frac{Z(\mathcal{A}, q_1)}{Z(\mathcal{A})} \prod_{n=1}^{N-1} \frac{\omega\left(q_n \xrightarrow{a/\cdot} q_{n+1}\right) Z(\mathcal{A}, q_{n+1})}{Z(\mathcal{A}, q_n)} \frac{\rho(q_N)}{Z(\mathcal{A}, q_N)}. \quad (4.48)$$

Notice that the state-specific allsums of all the inner states of the path (all states apart from q_1 and q_N) *cancel out* as the product moves over the transitions of the path. Additionally, the terms $Z(\mathcal{A}, q_1)$ and $Z(\mathcal{A}, q_N)$ cancel out with the definitions of $\lambda_{\mathcal{A}_L}$ and $\rho_{\mathcal{A}_L}$. This leaves us with

$$p_{\mathcal{A}_L}(\boldsymbol{\pi}) = \lambda(q_1) \frac{1}{Z(\mathcal{A})} \prod_{n=1}^{N-1} \omega\left(q_n \xrightarrow{a/\cdot} q_{n+1}\right) \rho(q_N) = p_{\mathcal{A}}(\boldsymbol{\pi}), \quad (4.49)$$

finishing the proof. ■

While Theorem 2.4.2 shows that any language model can be locally normalized, Theorem 4.1.1 shows that in the context of finite-state language models, the locally normalized version of a globally-normalized model is *also* a finite-state model.

Defining a Parametrized Globally Normalized Language Model

Having learned how an arbitrary normalizable finite-state language model can be normalized, we now discuss how models in this framework can be *parametrized* to enable fitting them to some training data. Crucial for parameterizing a globally normalized model is a **score function** $f^\delta_\theta : Q \times \Sigma \times Q \rightarrow \mathbb{R}$, which parametrizes the transitions between the states and thus determines the weights of the (accepting) paths. Additionally, we also parameterized the initial and final functions f^λ_θ and f^ρ_θ . These parametrized functions then define the automaton $\mathcal{A}_\theta \stackrel{\text{def}}{=} (\Sigma, Q, \delta_\theta, \lambda_\theta, \rho_\theta)$, where $\delta_\theta \stackrel{\text{def}}{=} \left\{ q_1 \xrightarrow{y/f^\delta_\theta(q_1, y, q_2)} q_2 \right\}$, $\lambda_\theta(q_i) \stackrel{\text{def}}{=} f^\lambda_\theta(q_i)$, and $\rho_\theta(q_\varphi) \stackrel{\text{def}}{=} f^\rho_\theta(q_\varphi)$. Note that we can parametrize the function f_θ in any way we want; for example, the function could be a neural network using distributed representations (we will see a similar example at the end of this section), or it could simply be a lookup table of weights. The fact that the function $f_\theta : (q_1, y, q_2)$ can only “look at” the identities of the states and the symbol might seem limiting; however, the states alone can encode a lot of information: for example, in n -gram models we describe below, they will encode the information about the previous $n - 1$ symbols and the transitions will then encode the probabilities of transitioning between such sequences of symbols.

The globally parametrized model then simply takes in any string $\mathbf{y} \in \Sigma^*$ and computes its stringsum value under the parametrized automaton, which in turn, as per Eq. (4.15), defines probabilities of the strings. The quantity $Z(\mathcal{A}_\theta)$ can be computed with the allsum algorithm discussed in §4.1.3. Importantly, since the algorithms for computing the string probabilities are differentiable, the model defined this way can also be *trained* with gradient-based learning as described in §3.2.3.

You might notice that this formulation does not exactly match the formulation of globally normalized models from Definition 2.4.2—the function $\mathcal{A} : \Sigma^* \rightarrow \mathbb{R}$ does not exactly match the form of an energy function as its values are not exponentiated as in Eq. (2.11). However, we tie this back to the definition of globally normalized models by defining an actual energy function as a simple *transformation* of the stringsum given by \mathcal{A}_θ . We can define the globally normalizing energy function $\hat{p}_{\text{GN}, \mathcal{A}_\theta}$ as

$$\hat{p}_{\text{GN}, \mathcal{A}_\theta}(\mathbf{y}) \stackrel{\text{def}}{=} -\log(\mathcal{A}(\mathbf{y})), \quad (4.50)$$

which can be easily seen to, after exponentiating it as in Eq. (2.11), result in the same expression as Eq. (4.15). With this, we have formulated finite-state language models as general globally normalized models.

Having introduced WFSAs as a formal and abstract computational model which can define a set of weighted strings, we now show how it can be used to explicitly model a particularly simple family of languages. We arrive at this family of language models when we impose a specific assumption on the set of *conditional distributions* of the language models that ensures that they are finite-state: the n -gram assumption.

4.1.4 Tightness of Finite-state Models

Any normalizable globally normalized finite-state language model is tight by definition because the sum of the scores over all *finite* strings is finite, and since they are normalized, they sum to 1. We, therefore, focus on *locally* normalized finite-state models and provide necessary and sufficient conditions for their tightness. Locally normalized finite-state models are exactly probabilistic WFSAs (Definition 4.1.17). Luckily, the tightness of probabilistic WFSAs can be easily characterized, as the following theorem shows.

Theorem 4.1.2. *A probabilistic WFSA is tight if and only if all accessible states are also co-accessible.*

Proof. We prove each direction in turn.

(\Rightarrow): Assume the WFSA is tight. Let $q \in Q$ be an accessible state, which means q can be reached after a finite number of steps with positive probability. By tightness assumption, then there must be a positive probability path from q to termination, or else the WFSA will not be able to terminate after reaching q , resulting in non-tightness. This means that q is also co-accessible. So, assuming that the WFSA is tight, every accessible state is also co-accessible.

(\Leftarrow): Assume that all accessible states are co-accessible. First, one may consider a Markov chain consisting only of the set of accessible states $Q_A \subseteq Q$, since all other states will have probability 0 at every step. Recall a fundamental result in finite-state Markov chain theory which states that, if there exists a unique absorbing state which is reachable from every state, then the Markov process is absorbed by this state with probability 1 (see, e.g., Theorem 11.3 in Grinstead and Snell, 1997). We already have that

- EOS is an absorbing state, and that
- by assumption, every state in Q_A is co-accessible which implies that they can reach EOS.

Hence, it remains to show that EOS is the *unique* absorbing state. Suppose there is another state (or group of states) in Q_A distinct from EOS that is absorbing, i.e., cannot leave once entered. Then, these states cannot reach EOS by assumption, which means they are not co-accessible, contradicting the assumption that every state in Q_A is co-accessible. Hence, EOS is the only absorbing state in Q_A and by the property of an absorbing Markov chain, the process is absorbed by EOS with probability 1. In other words, the WFSFA is tight. ■

Notice that trimming a PFSA results in a model that satisfies $\rho(q) + \sum_{q \xrightarrow{a/w} q'} w \leq 1$, but might no longer achieve equality as required by Definition 4.1.17. We call such models substochastic WFSAs.

Definition 4.1.23 (Substochastic Weighted Finite-State Automaton). *A WFSFA $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ is **substochastic** if for all $q \in Q$ and all outgoing transitions $q \xrightarrow{a/w} q' \in \delta$ it holds that*

$$\lambda(q) \geq 0 \quad (4.51)$$

$$\rho(q) \geq 0 \quad (4.52)$$

$$w \geq 0 \quad (4.53)$$

and

$$\rho(q) + \sum_{q \xrightarrow{a/w} q'} w \leq 1. \quad (4.54)$$

We can then express the termination probability of a WFSFA in simple linear algebra terms.

Theorem 4.1.3. *Let \mathbf{T}' be the transition sum matrix of a trimmed substochastic WFSFA. Then $\mathbf{I} - \mathbf{T}'$ is invertible and $p(\mathbf{x} \in \Sigma^*) = \vec{\lambda}'^\top (\mathbf{I} - \mathbf{T}')^{-1} \vec{\rho}' \leq 1$.*

In the following, we will make use of the spectral radius of a matrix.

Definition 4.1.24 (Spectral radius). *The **spectral radius** of a matrix $\mathbf{M} \in \mathbb{C}^{N \times N}$ with eigenvalues $\lambda_1, \dots, \lambda_N$ is defined as*

$$\rho_s(\mathbf{M}) \stackrel{\text{def}}{=} \max\{|\lambda_1|, \dots, |\lambda_N|\}. \quad (4.55)$$

To prove Theorem 4.1.3, we will make use of the following useful lemma.

Lemma 4.1.3. *Let \mathbf{T}' be the transition sum matrix of a trimmed substochastic WFSFA, then $\rho_s(\mathbf{T}') < 1$.*

Proof. To begin with, we wish to apply the following result which connects the row sums of a matrix to its spectral radius. Below, \mathcal{M}_N denotes the set of $N \times N$ matrices, and $\|\mathbf{A}\|_\infty = \max_{1 \leq n \leq N} \sum_{i=1}^N |\mathbf{A}_{ni}|$ denotes the infinity matrix norm.

Proposition 4.1.1 (§6.2.P8; Horn and Johnson, 2012). *For any $\mathbf{A} \in \mathcal{M}_N$, $\rho_s(\mathbf{A}) \leq \|\mathbf{A}\|_\infty$. Additionally, if \mathbf{A} is irreducible and not all absolute row sums of \mathbf{A} are equal, then $\rho_s(\mathbf{A}) < \|\mathbf{A}\|_\infty$.*

However, the transition sum matrix \mathbf{P} of a substochastic WFSFA may be reducible whereas the irreducibility condition in Proposition 4.1.1 cannot be dropped. Hence, we need to “decompose” \mathbf{T}' in a way to recover irreducibility. We use the *Frobenius normal form* (also known as *irreducible normal form*) to achieve this.

Proposition 4.1.2 (§8.3.P8; Horn and Johnson, 2012). *Let $\mathbf{A} \in \mathcal{M}_N$ be non-negative. Then, either \mathbf{A} is irreducible or there exists a permutation matrix \mathbf{P} such that*

$$\mathbf{P}^\top \mathbf{A} \mathbf{P} = \begin{bmatrix} \mathbf{A}_1 & & * \\ & \ddots & \\ \mathbf{0} & & \mathbf{A}_K \end{bmatrix} \quad (4.56)$$

*is block upper triangular, and each diagonal block is irreducible (possibly a 1-by-1 zero matrix). This is called an **Frobenius normal form** (or **irreducible normal form**) of A . Additionally, $\Lambda(A) = \Lambda(\mathbf{A}_1) \cup \dots \cup \Lambda(\mathbf{A}_K)$ where $\Lambda(\cdot)$ denotes the set of eigenvalues of a matrix.*

Notice that, by way of a similarity transformation via a permutation matrix, the Frobenius normal form is equivalent to a relabeling of the states in the trimmed WFSFA in the sense of

$$(\mathbf{P}^\top \vec{\lambda}')^\top (\mathbf{P}^\top \mathbf{T}' \mathbf{P})^K (\mathbf{P}^\top \vec{\rho}') = (\vec{\lambda}'^\top \mathbf{P}) (\mathbf{P}^\top \mathbf{T}'^K \mathbf{P}) (\mathbf{P}^\top \vec{\rho}') \quad (4.57a)$$

$$= \vec{\lambda}'^\top \mathbf{T}'^K \vec{\rho}' \quad (4.57b)$$

where the equalities follow from the fact that the inverse of a permutation matrix \mathbf{P} is its transpose. Hence, with an appropriate relabeling, we may assume without loss of generality that \mathbf{P} is already put into a Frobenius normal form

$$\mathbf{T}' = \begin{bmatrix} \mathbf{T}'_1 & & * \\ & \ddots & \\ \mathbf{0} & & \mathbf{T}'_K \end{bmatrix} \quad (4.58)$$

where each \mathbf{T}'_k is irreducible.

Since the transition sum matrix \mathbf{T}' of a trimmed substochastic WFSFA is a substochastic matrix, each \mathbf{T}'_k is also substochastic. In fact, each \mathbf{T}'_k is *strictly* substochastic, meaning that there is at least a row that sums to less than 1. To see this, suppose to the contrary that there is a probabilistic \mathbf{T}'_k . Since the WFSFA is trimmed, every state is both accessible and co-accessible. Being accessible implies that there is a positive probability of reaching every state in \mathbf{T}'_k . However, the probabilisticity of \mathbf{T}'_k forces the corresponding $\vec{\rho}'$ entries to be 0. Hence, none of these states can transition to EOS, meaning that they're not co-accessible, contradicting the assumption. Hence, every \mathbf{T}'_k is strictly substochastic and has at least one strictly less than 1 row sum. Then, either all row sums of \mathbf{T}'_k are less than 1 or some row sums are 1 and some are less than 1. In either cases, Proposition 4.1.1 implies that $\rho_s(\mathbf{T}'_k) < 1$ for all $1 \leq k \leq K$. Finally, as Proposition 4.1.2 entails, $\rho_s(\mathbf{T}') = \max\{\rho_s(\mathbf{T}'_1), \dots, \rho_s(\mathbf{T}'_K)\}$ where each $\rho_s(\mathbf{T}'_k) < 1$. Hence, $\rho_s(\mathbf{T}') < 1$. ■

We now use the stated results to finally prove Theorem 4.1.3.

Proof. By Lemma 4.1.3, $\rho_s(\mathbf{T}') < 1$, in which case $\mathbf{I} - \mathbf{T}'$ is invertible and the Neumann series $\mathbf{I} + \mathbf{T}' + \mathbf{T}'^2 + \dots$ converges to $(\mathbf{I} - \mathbf{T}')^{-1}$ (§5.6, Horn and Johnson, 2012). Hence, we can write $(\mathbf{I} - \mathbf{T}')^{-1} = \sum_{k=0}^{\infty} \mathbf{T}'^k$. Then,

$$p(\Sigma^*) = \sum_{k=0}^{\infty} P(\Sigma^k) \quad (4.59a)$$

$$= \sum_{k=0}^{\infty} \vec{\lambda}'^\top \mathbf{T}'^k \vec{\rho}' \quad (4.59b)$$

$$= \vec{\lambda}'^\top \left(\sum_{k=0}^{\infty} \mathbf{T}'^k \right) \vec{\rho}' \quad (4.59c)$$

$$= \vec{\lambda}'^\top (\mathbf{I} - \mathbf{T}')^{-1} \vec{\rho}'. \quad (4.59d)$$

■

4.1.5 The n -gram Assumption and Subregularity

We now turn our attention to one of the first historically significant language modeling frameworks: n -gram models. While they are often taught completely separately from (weighted) finite-state automata, we will see shortly that they are simply a special case of finite-state language models and thus all results for the more general finite-state language models also apply to the specific n -gram models as well.

As we saw in Theorem 2.4.2, we can factorize the language model p_{LM} for $\mathbf{y} = y_1 \dots y_T \in \Sigma^*$ as

$$p_{\text{LM}}(\mathbf{y}) = p_{\text{LN}}(\mathbf{y}) = p_{\text{SM}}(\text{EOS} \mid \mathbf{y}) \prod_{t=1}^T p_{\text{SM}}(y_t \mid \mathbf{y}_{<t}), \quad (4.60)$$

where $p_{\text{SM}}(y \mid \mathbf{y})$ are specified by a locally normalized model (Definition 2.4.5).

Recall that SMs specify individual conditional distributions of the next symbol y_t given the previous $t - 1$ symbols for *all possible* t . However, as t grows and the history of seen tokens accumulates, the space of possible histories (sequences of strings to condition on) grows very large (and indeed infinite as $t \rightarrow \infty$). This makes the task of modeling individual conditional distributions for large t computationally infeasible. One way to make the task more manageable is by using the n -gram assumption.

Assumption 4.1.1 (n -gram assumption). *In words, the **n -gram assumption** states that the probability of a word y_t only depends on $n - 1$ previous words $y_{t-1}, \dots, y_{t-n+1}$ where $y_0 \stackrel{\text{def}}{=} \text{BOS}$. Notationally, we can write the n -gram assumption as a conditional independence assumption i.e.,*

$$p_{\text{SM}}(y_t \mid \mathbf{y}_{<t}) \stackrel{\text{def}}{=} p_{\text{SM}}(y_t \mid y_{t-1} \dots y_{t-n+1}) \stackrel{\text{def}}{=} p_{\text{SM}}(y_t \mid \mathbf{y}_{t-n-1:t-1}) \quad (4.61)$$

The sequence $y_{t-1} \dots y_{t-n+1}$ is often called the **history** or the **context**.

In plain English, this means that the probability of a token only depends on the previous $n - 1$ tokens. n -gram assumption is, therefore, an alias of $(n - 1)^{\text{th}}$ order Markov assumption in the language modeling context.

Handling edge cases by padding. Given our definition in Eq. (4.61) where the conditional probability $p_{\text{SM}}(y_t \mid \mathbf{y}_{t-n-1:t-1})$ depends on exactly $n - 1$ previous symbols, we could run into an issue with negative indices for $t < n$. To handle edge cases for $t < n$, we will **pad** the sequences with the BOS symbols at the beginning, that is, we will assume that the sequences $\bar{y}_1 \dots \bar{y}_t$ for $t < n - 1$ are “transformed” as

$$\bar{y}_1 \bar{y}_2 \dots \bar{y}_t \mapsto \underbrace{\text{BOS} \dots \text{BOS}}_{n-1-t \text{ times}} \bar{y}_1 \bar{y}_2 \dots \bar{y}_t \quad (4.62)$$

Notice that with such a transformation, we always end up with strings of length $n - 1$, which is exactly what we need for conditioning in an n -gram model. In the following, we will assume that all such sequences are already transformed, but at the same time, we will assume that

$$p_{\text{SM}} \left(\bar{y} \mid \underbrace{\text{BOS} \dots \text{BOS}}_{n-1-t \text{ times}} \bar{y}_1 \bar{y}_2 \dots \bar{y}_t \right) = \bar{y}_0 \bar{y}_1 \bar{y}_2 \dots \bar{y}_t \quad (4.63)$$

By definition, n -gram language models can only model dependencies spanning n tokens or less. By limiting the length of the relevant context when determining $p_{\text{SM}}(y_t \mid \mathbf{y}_{<t})$ to the previous n tokens, the n -gram assumption limits the number of possible probability *distributions* that need to be tracked to $\mathcal{O}(|\Sigma|^{n-1})$.

A particularly simple case of the n -gram model is the **bigram** model where $n = 2$, which means that the probability of the next word only depends on the previous one, i.e., $p_{\text{SM}}(y_t \mid \mathbf{y}_{<t}) = p_{\text{SM}}(y_t \mid y_{t-1})$.¹⁷

Example 4.1.5 (A simple bigram model). *Let us look at a specific example of a simple bigram model. Suppose our vocabulary consists of the words “large”, “language”, and “models”, thus, $|\Sigma| = 3$. To specify the bigram model, we have to define the conditional probabilities $p_M(y_j \mid y_i)$ for $y_i \in \Sigma \cup \{\text{BOS}, \text{EOS}\}$ and*

¹⁷What would the uni-gram ($n = 1$) model look like? What conditional dependencies between words in a sentence could be captured by it?

$y_j \in \Sigma \cup \{\text{EOS}\}$ (remember that we do not have to model the probability of the next token being BOS). In the case of bigrams, we can represent those in a table, where the entry at position i, j represents the probability $p_M(y_j \mid y_i)$:

	"large"	"language"	"models"	"EOS"
BOS	0.4	0.2	0.2	0.2
"large"	0.1	0.4	0.2	0.3
"language"	0.1	0.1	0.4	0.4
"models"	0.2	0.2	0.1	0.5
"EOS"	0.4	0.2	0.2	0.2

Under our model, the probability of the sentence "large language models" would be

$$\begin{aligned}
 & p_{SM}(\text{"large"} \mid \text{BOS}) \\
 & \cdot p_{SM}(\text{"language"} \mid \text{"large"}) \\
 & \cdot p_{SM}(\text{"models"} \mid \text{"language"}) \\
 & \cdot p_{SM}(\text{EOS} \mid \text{"models"}) \\
 & = 0.4 \cdot 0.4 \cdot 0.4 \cdot 0.5 = 0.032
 \end{aligned}$$

while the probability of the sentence "large large large" would be

$$\begin{aligned}
 & p_{SM}(\text{"large"} \mid \text{BOS}) \\
 & \cdot p_{SM}(\text{"large"} \mid \text{"large"}) \\
 & \cdot p_{SM}(\text{"large"} \mid \text{"large"}) \\
 & \cdot p_{SM}(\text{EOS} \mid \text{"large"}) \\
 & = 0.4 \cdot 0.1 \cdot 0.1 \cdot 0.3 = 0.0012.
 \end{aligned}$$

Note that the probabilities in the above table are made up and not completely reasonable. A real n -gram model would not allow for probabilities of exactly 0 to avoid pathological behavior.

Representing n -gram Models as WFSA

We define n -gram language models as models which only consider a finite amount of context when defining the conditional probabilities of the next token. This means that the set of possible conditional distributions $p_{SM}(y \mid \mathbf{y})$ is also finite which very naturally connects them to weighted finite-state automata—indeed, every n -gram language model is a WFSA—specifically, a probabilistic finite-state automaton (or a substochastic one). We will make this connection more formal in this subsection, thus formally showing that n -gram models are indeed finite-state. Note that this is different from §4.1.3, where we discussed how to parametrize a general WFSA and use it as a globally normalized model—in contrast, in this section, we consider how to fit a (locally normalized) n -gram model into the finite-state framework.

The intuition behind the connection is simple: the finite length of the context implies a finite number of histories we have to model. These histories represent the different states the corresponding automaton can reside in at any point. Given any history \mathbf{y} with $|\mathbf{y}| < n$ and the state $q \in Q$ representing \mathbf{y} , then, the conditional distribution of the next token given \mathbf{y} dictate the transition weights into the next states in the WFSA, representing the new, updated history of the input.

Importantly, since we want PFSA to represent globally-normalized models, we will also remove the EOS symbol from the n -gram model before transforming it into a PFSA—as the remark above about the relationship between the EOS symbol and the final states hints, the latter will fill in the role of the EOS symbol. The way we do that is the following. From the semantics of the EOS symbol discussed in the section on tightness (cf. Eq. (2.44)), we also know that to model the probability distribution over finite strings in Σ^* , we only require to keep track of strings up to the first occurrence of the EOS symbol. Therefore, when converting a given n -gram model to a WFSA, we will only model sequences up to the first occurrence of the special symbol, meaning that EOS will never occur in the *context* of any conditional distribution $p_{SM}(\bar{y} \mid \mathbf{y})$. We now detail this construction.

Let p_{LN} be a well-defined n -gram language model specified by conditional distributions p_{SM} as defined by Assumption 4.1.1. We will now construct a WFSA representing p_{LN} . Intuitively, its states will represent all possible sequences of words of length n while the transitions between the states q_1 and q_2 will correspond to the possible transitions between the n -grams which those represent. This means that the only possible (positively weighted) transitions will be between the n -grams which can follow each other, i.e. $\mathbf{y}_{t-n:t-1}$ and $\mathbf{y}_{t-n+2:t}$ for some $y_{t-n}, y_t \in \bar{\Sigma}$ (until the first occurrence of EOS). The transition's weight will depend on the probability of observing the “new” word y_0 in the second n -gram given the starting n -gram $y_{-n}y_{-(n-1)} \dots y_{-1}$. Further, the final weights of the states will correspond to *ending* the string in them. In p_{LN} , this is modeled as the probability of observing EOS given the context $\mathbf{y}_{t-n:t-1}$ —this, therefore, is set as the final weight of the state representing the history $\mathbf{y}_{t-n:t-1}$. Formally, we can map a n -gram model into a WFSA $\mathcal{A} = (\Sigma_{\mathcal{A}}, Q_{\mathcal{A}}, \delta_{\mathcal{A}}, \lambda_{\mathcal{A}}, \rho_{\mathcal{A}})$ by constructing \mathcal{A} as follows.

- Automaton's alphabet:

$$\Sigma_{\mathcal{A}} \stackrel{\text{def}}{=} \Sigma \quad (4.64)$$

- The set of states:

$$Q_{\mathcal{A}} \stackrel{\text{def}}{=} \bigcup_{t=0}^{n-1} \{\text{BOS}\}^{n-1-t} \times \Sigma^t \quad (4.65)$$

- The transitions set

$$\delta_{\mathcal{A}} \stackrel{\text{def}}{=} \{ \mathbf{y}_{t-n:t-1} \xrightarrow{y_t / p_{\text{SM}}(y_t | \mathbf{y}_{t-n:t-1})} \mathbf{y}_{t-n+1:t} \mid \mathbf{y}_{t-n+1:t-1} \in \bigcup_{t=0}^{n-2} \{\text{BOS}\}^{n-2-t} \times \Sigma^t; y_{t-n}, y_t \in \Sigma \} \quad (4.66)$$

- The initial function:

$$\lambda_{\mathcal{A}} : \mathbf{y} \mapsto \begin{cases} 1 & \text{if } \mathbf{y} = \underbrace{\text{BOS} \dots \text{BOS}}_{n-1 \text{ times}} \\ 0 & \text{otherwise} \end{cases} \quad (4.67)$$

- The final function

$$\rho_{\mathcal{A}} : \mathbf{y} \mapsto p_{\text{SM}}(\text{EOS} \mid \mathbf{y}), \mathbf{y} \in Q_{\mathcal{A}} \quad (4.68)$$

The definition of the states set $Q_{\mathcal{A}}$ captures exactly the notion of padding with the BOS symbol for handling the edge cases we described above. This shows that n -gram language models are indeed finite-state (we leave the formal proof showing that $L(\mathcal{A}) = L(p_{\text{LN}})$ to the reader).

Defining a n -gram language model through a parametrized WFSA. We now consider how we can use the framework of WFSA to define a more “flexible” *parametrized globally* normalized model. In this case, we do not start from an existing locally normalized set of distributions forming p_{SM} . Rather, we would like to model the “suitability” of different n -grams following each other—that is, we would like to somehow *parametrize* the probability that some n -gram \mathbf{y}' will follow an n -gram \mathbf{y} without having to worry about normalizing the model at every step. This will allow us to then fit the probability distributions of the model to those in the data, e.g., with techniques described in §3.2.3. Luckily, the flexibility of the WFSA modeling framework allows us to do exactly that.

Subregularity

We saw that language models implementing the very natural n -gram assumption can be represented using weighted finite-state automata. However, n -gram models do not “need the full expressive power” of WFSA—they can actually be modeled using even simpler machines than finite-state automata. This, along with several other examples of simple families of formal languages, motivates the definition of *subregular* languages.

Definition 4.1.25. A language is **subregular** if it can be recognized by a finite-state automaton or any weaker machine.

Most subregular languages can indeed be recognized by formalisms which are much simpler than FSAs. Many useful and interesting classes of subregular languages have been identified—recently, especially in the field of phonology. Naturally, due to their simpler structure, they also allow for more efficient algorithms—this is why we always strive to represent a language with the simplest formalism that still captures it adequately. See Jäger and Rogers (2012); Avcu et al. (2017) for comprehensive overviews of subregular languages.

Subregular languages actually form multiple hierarchies of complexity within regular languages. Interestingly, n -gram models fall into the simplest level of complexity in one of the hierarchies, directly above *finite* languages. This class of subregular languages is characterized by patterns that depend solely on the blocks of symbols that occur *consecutively* in the string, which each of the blocks considered independently of the others—it is easy to see that n -gram models intuitively fall within such languages. This family of subregular languages is suggestively called *strictly local languages*.

Definition 4.1.26 (Strictly local languages). A language L is **strictly n -local** (SL_n) if, for every string \mathbf{y} of length $|\mathbf{y}| = n - 1$, and all strings $\mathbf{x}_1, \mathbf{x}_2, \mathbf{z}_1, \mathbf{z}_2 \in \Sigma^*$, it holds that if $\mathbf{x}_1\mathbf{y}\mathbf{z}_1 \in L$ and $\mathbf{x}_2\mathbf{y}\mathbf{z}_2 \in L$, then also $\mathbf{x}_1\mathbf{y}\mathbf{z}_2 \in L$ (and $\mathbf{x}_2\mathbf{y}\mathbf{z}_1 \in L$).

A language is **strictly local** (SL) if it is strictly n -local for any n .

Note that we could of course also define this over with the EOS-augmented alphabet $\bar{\Sigma}$. You can very intuitively think of this definition as postulating that the history more than n symbols back does not matter anymore for determining or specifying whether a string is in a language (or its weight, in the weighted case)—this is exactly what the n -gram assumption states.

4.1.6 Representation-based n -gram Models

So far, we have mostly talked about the conditional probabilities and the WFSA weights defining a language model very abstractly. Apart from describing how one can generally parametrize the weights of the underlying WFSA with the scoring function in §4.1.5, we only discussed what values the weights can take for the language model to be well-defined and what implications that has on the distribution defined by the WFSA. In this section, we consider for the first time what an actual implementation of a finite-state, or more precisely, a n -gram language model might look like. Concretely, we will define our first parameterized language model in our General language modeling framework (cf. §3.1) by defining a particular form of the encoding function enc as a simple multi-layer feed-forward neural network.¹⁸

However, before we dive into that, let us consider as an alternative possibly the simplest way to define a (locally normalized) n -gram language model: by directly parametrizing the probabilities of each of the symbols y in the distribution $p_{\text{SM}}(\bar{y} \mid \bar{\mathbf{y}})$ for any context $\bar{\mathbf{y}}$, that is

$$\boldsymbol{\theta} \stackrel{\text{def}}{=} \left\{ \theta_{y|\mathbf{y}} \stackrel{\text{def}}{=} p_{\text{SM}}(y \mid \mathbf{y}) \mid y \in \bar{\Sigma}, \mathbf{y} \in \bar{\Sigma}^{n-1}, \theta_{y|\mathbf{y}} \geq 0, \sum_{y' \in \bar{\Sigma}} \theta_{y'|\mathbf{y}} = 1 \right\}. \quad (4.69)$$

The following proposition shows that the maximum likelihood solution (Eq. (3.60)) to this parametrization is what you would probably expect.

Proposition 4.1.3. The MLE solution of Eq. (4.69) is

$$p_{\text{SM}}(y_n \mid \mathbf{y}_{<n}) = \frac{C(y_1, \dots, y_n)}{C(y_1, \dots, y_{n-1})} \quad (4.70)$$

whenever the denominator > 0 , where $C(y_1, \dots, y_n)$ denotes the number of occurrences of all possible strings of the form y_1, \dots, y_n and $C(y_1, \dots, y_{n-1})$ denotes the number of occurrences of all possible strings of the form y_1, \dots, y_{n-1} .

¹⁸While we introduce particular architectures of neural networks, for example, recurrent neural networks and transformers later in Chapter 4, we assume some familiarity with neural networks in general. See Chapter 6 of Goodfellow et al. (2016) for an introduction.

Proof. Let $\mathcal{D} \stackrel{\text{def}}{=} \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(M)}\}$ be the training dataset. The log-likelihood of a single example $\mathbf{y}^{(m)}$ is

$$\log(p_{\text{LN}}(\mathbf{y})) = \log \left(\prod_{t=1}^{|\mathbf{y}^{(m)}|} p_{\text{SM}} \left(y_t^{(m)} \mid y_{t-n:t-1}^{(m)} \right) \right) \quad (4.71)$$

$$= \sum_{t=1}^{|\mathbf{y}^{(m)}|} \log p_{\text{SM}} \left(y_t^{(m)} \mid y_{t-n:t-1}^{(m)} \right) \quad (4.72)$$

which means that the log-likelihood of the entire dataset is

$$\ell\ell(\mathcal{D}) = \sum_{m=1}^M \sum_{t=1}^{|\mathbf{y}^{(m)}|} \log p_{\text{SM}} \left(y_t^{(m)} \mid y_{t-n:t-1}^{(m)} \right) \quad (4.73)$$

$$= \sum_{m=1}^M \sum_{t=1}^{|\mathbf{y}^{(m)}|} \log \theta_{y_n | \mathbf{y}^{(m)}_{<n}}. \quad (4.74)$$

Exercise 4.2 asks you to show that this can be rewritten with the **token to type switch** as

$$\ell\ell(\mathcal{D}) = \sum_{\substack{\mathbf{y} \\ |\mathbf{y}|=n}} C(\mathbf{y}) \theta_{y_n | \mathbf{y}^{<n}}. \quad (4.75)$$

The maximum likelihood parameters can then be determined using Karush–Kuhn–Tucker (KKT) conditions¹⁹ to take into account the non-negativity and local normalization constraints:

$$\nabla_{\boldsymbol{\theta}} \left(\ell\ell(\mathcal{D}) - \sum_{\substack{\mathbf{y} \in \Sigma^* \\ |\mathbf{y}|=n-1}} \lambda_{\mathbf{y}} \left(\sum_{y \in \Sigma} \theta_{y|\mathbf{y}} - 1 \right) - \sum_{\substack{\mathbf{y} \in \Sigma^* \\ |\mathbf{y}|=n-1}} \eta_{\mathbf{y}} \theta_{y|\mathbf{y}} \right) = 0. \quad (4.76)$$

Recall that the KKT conditions state that a $\boldsymbol{\theta}$ is an optimal solution of $\ell\ell$ if and only if $(\boldsymbol{\theta}, \{\lambda_{\mathbf{y}}\}_{\mathbf{y} \in \Sigma^{n-1}}, \{\eta_{\mathbf{y}}\}_{\mathbf{y} \in \Sigma^{n-1}})$ satisfy Eq. (4.76). Since this is simply a sum over the dataset with no interactions of parameters for individual contexts \mathbf{y} with $|\mathbf{y}| = n-1$ in $\theta_{y|\mathbf{y}}$, it can be solved for each context \mathbf{y} individually.

Moreover, as you are asked to show Exercise 4.3, it holds that

$$\sum_{y' \in \Sigma} C(y_1 \dots y_{n-1} y') = C(y_1 \dots y_{n-1}) \quad (4.77)$$

for any $\mathbf{y} = y_1 \dots y_{n-1} \in \Sigma^{n-1}$. This leaves us with the following system for each $\mathbf{y} \in \Sigma^{n-1}$:

$$\sum_{y' \in \Sigma} C(\mathbf{y} y') \log \theta_{y'|\mathbf{y}} - \lambda_{\mathbf{y}} \left(\sum_{y' \in \Sigma} \theta_{y'|\mathbf{y}} - 1 \right) - \eta_{\mathbf{y}} \theta_{y|\mathbf{y}} = 0. \quad (4.78)$$

It is easy to confirm that $\theta_{y|\mathbf{y}} = \frac{C(\mathbf{y} y)}{C(\mathbf{y})}$ with $\lambda_{\mathbf{y}} = C(\mathbf{y})$ and $\eta_{\mathbf{y}} = 0$ is a saddle point of Eq. (4.76). This means that $\theta_{y|\mathbf{y}} = \frac{C(\mathbf{y} y)}{C(\mathbf{y})}$ is indeed the maximum likelihood solution. ■

This results in a locally normalized n -gram model. To avoid issues with division-by-zero and assigning 0 probability to unseen sentences, we can employ methods such as *smoothing* and *backoff*, which are beyond the scope of the course.²⁰

While this model might seem like an obvious choice, it comes with numerous drawbacks. To see what can go wrong, consider the following example.

¹⁹See https://en.wikipedia.org/wiki/Karush%E2%80%93Kuhn%E2%80%93Tucker_conditions.

²⁰See (Chen and Goodman, 1996) and Chapter 4 in (Jurafsky and Martin, 2009).

Example 4.1.6. Suppose we have a large training corpus of sentences, among which sentences like “We are going to the shelter to adopt a dog.”, “We are going to the shelter to adopt a puppy.”, and “We are going to the shelter to adopt a kitten.”, however, without the sentence “We are going to the shelter to adopt a cat.” Fitting an n -gram model using the count statistics and individual tables of conditional probabilities $p_{SM}(y \mid \mathbf{y})$, we would assign the probability

$$p_{SM}(y_t = \text{cat} \mid \mathbf{y}_{<t} = \text{We are going to the shelter to adopt a})$$

the value 0 (or some “default” probability if we are using smoothing). However, the words “dog”, “puppy”, “kitten”, and “cat” are semantically very similar—they all describe pets often found in shelters. It would therefore be safe to assume that the word “cat” is similarly probable given the context “We are going to the shelter to adopt a” as the other three words observed in the training dataset. However, if we estimate all the conditional probabilities independently, we have no way of using this information—the words have no relationship in the alphabet, they are simply different indices in a lookup table. Additionally, statistics gathered for the sentences above will not help us much when encountering very similar sentences, such as “We went to a nearby shelter and adopted a kitten.” The issue is that there are simply many ways of expressing similar intentions. We would thus like our language models to be able to generalize across different surface forms and make use of more “semantic” content of the sentences and words. However, if the model is parametrized as defined in Eq. (4.69), it is not able to take advantage of any such relationships.

The model defined by Eq. (4.69) is therefore unable to take into account the relationships and similarities between words. The general modeling framework defined in §3.1 allows us to remedy this using the **distributed word representations**. Recall that, in that framework, we associate each word y with its vector representation $\mathbf{e}(y)$ (its *embedding*), and we combine those into the embedding matrix \mathbf{E} . Importantly, word embeddings are simply additional parameters of the model and can be *fit* on the training dataset *together* with the language modeling objective. One of the first successful applications of $\text{enc}(\cdot)$ is due to Bengio et al. (2003), which we discuss next.

To be able to use the embeddings in our general framework, we now just have to define the concrete form of the context-encoding function enc . In the case of the neural n -gram model which we consider here and as defined by (Bengio et al., 2003), the representations of the context $\mathbf{y}_{<t}$, $\text{enc}(\mathbf{y}_{<t})$, are defined as the output of a neural network which looks at the previous $n - 1$ words in the context:

$$\text{enc}(\mathbf{y}_{<t}) \stackrel{\text{def}}{=} \text{enc}(y_{t-1}, y_{t-2}, \dots, y_{t-n+1}), \quad (4.79)$$

where enc is a neural network we define in more detail shortly. The full language model is therefore defined through the conditional distributions

$$p_{SM}(\bar{y}_t \mid \bar{\mathbf{y}}_{<t}) \stackrel{\text{def}}{=} \text{softmax} \left(\text{enc}(\bar{y}_{t-1}, \bar{y}_{t-2}, \dots, \bar{y}_{t-n+1})^\top \mathbf{E} + \mathbf{b} \right)_{\bar{y}_t} \quad (4.80)$$

resulting in the locally normalized model

$$p_{LN}(\mathbf{y}) = \text{softmax} \left(\text{enc}(\bar{y}_T, \bar{y}_{T-1}, \dots, \bar{y}_{T-n+2})^\top \mathbf{E} + \mathbf{b} \right)_{\text{EOS}} \quad (4.81)$$

$$\cdot \prod_{t=1}^T \text{softmax} \left(\text{enc}(y_{t-1}, y_{t-2}, \dots, y_{t-n+1})^\top \mathbf{E} + \mathbf{b} \right)_{y_t} \quad (4.82)$$

for $\mathbf{y} \in \Sigma^*$.

Importantly, notice that although this is a *neural* model, it is nonetheless still an n -gram model with finite context—Eq. (4.79) is simply a restatement of the n -gram assumption in terms of the neural encoding function enc . It therefore still suffers from some of the limitations of regular n -gram models, such as the inability to model dependencies spanning more than n words. However, it solves the problems encountered in Example 4.1.6 by considering word similarities and *sharing* parameters across different contexts in the form of an encoding function rather than a lookup table.

While encoding function enc in Eq. (4.79) could in principle take any form, the original model defined in Bengio et al. (2003) defines the output as for the string $\mathbf{y} = y_t, y_{t-1}, \dots, y_{t-n+1}$ as

$$\text{enc}(y_t, y_{t-1}, \dots, y_{t-n+1}) \stackrel{\text{def}}{=} \mathbf{b} + \mathbf{W}\mathbf{x} + \mathbf{U} \tanh(\mathbf{d} + \mathbf{H}\mathbf{x}), \quad (4.83)$$

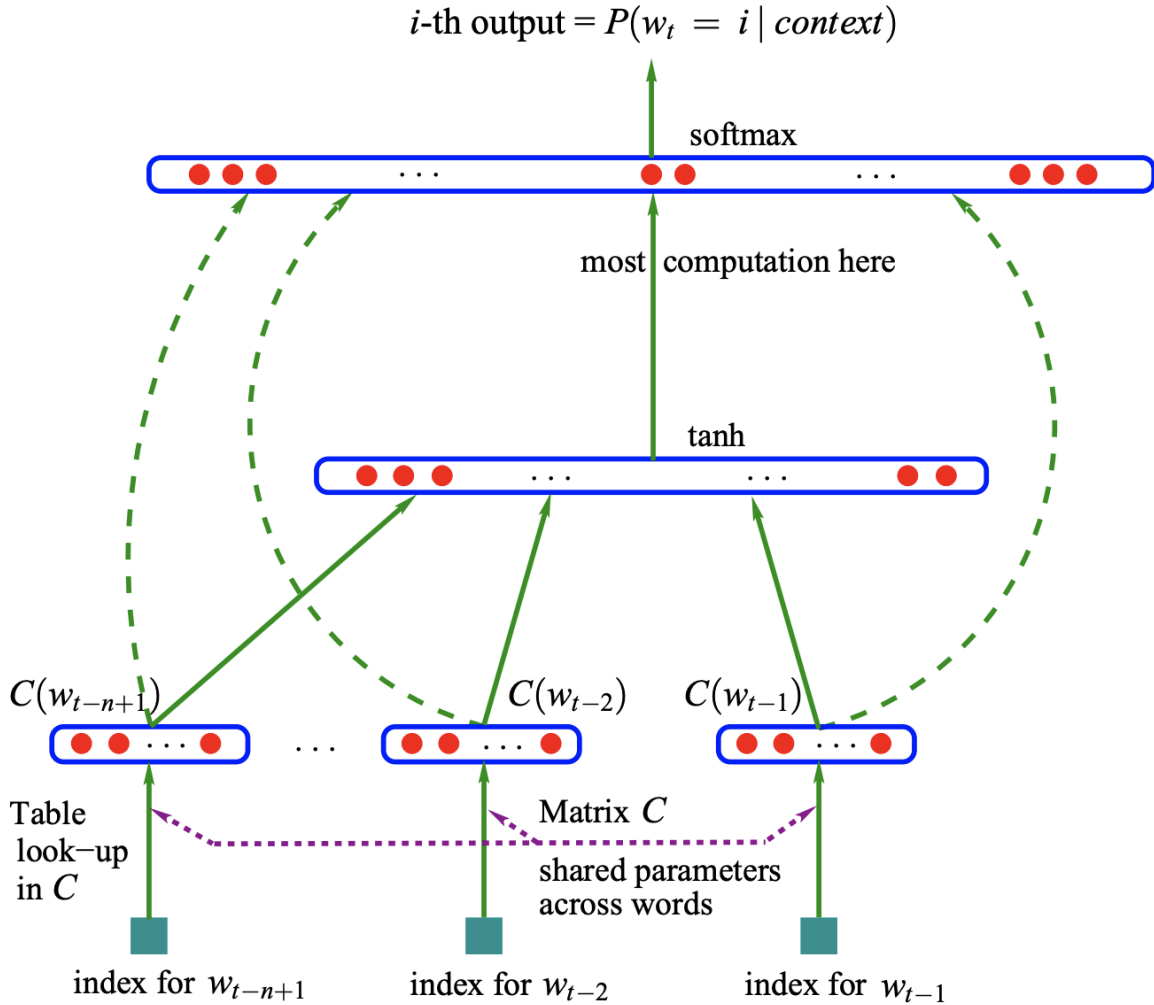


Figure 4.4: A pictorial depiction of the n -gram neural language model from the original publication (Bengio et al., 2003). Note that the quantity $C(w)$ corresponds to $\mathbf{e}(y)$ for a word y in our notation.

where $\mathbf{x} \stackrel{\text{def}}{=} \text{concat}(\mathbf{e}(y_t), \mathbf{e}(y_{t-1}), \dots, \mathbf{e}(y_{t-n+1}))$ denotes the concatenation of the context symbol embeddings into a long vector of size $(n-1) \cdot R$, and \mathbf{b} , \mathbf{d} , \mathbf{W} , and \mathbf{U} define the parameters of the encoding function. This completes our definition of the model in the general language modeling framework—the model can then simply be trained on the language modeling objective as defined in §3.2.2.

We can also see that such a model also reduces the number of parameters required to specify a n -gram model: whereas a lookup-table-based n -gram model with no parameter sharing requires $\mathcal{O}(|\Sigma|^n)$ parameters to be defined, the number of parameters required by a representation-based n -gram model scales *linearly* with n —all we have to do is add additional rows to the matrices defined in Eq. (4.83). We will later see how this can be reduced to a *constant* number of parameters w.r.t. the sequence length in the case of recurrent neural networks in §5.1.2.

Pictorially, we can imagine the model as depicted in Fig. 4.4 (taken from the original publication). This shows that the n -gram modeling framework is not limited to counting co-occurrence statistics. The model from Eq. (4.79) can also be represented by a WFSA just like the simpler models we discussed above, with the weights on the transitions parametrized by the neural network. This allows us to both understand well with insights from formal language theory, as well as to train it in a flexible way allowed for by the non-linear encoding function. However, the model from Eq. (4.79) is still limited to statistics of the last n tokens or less.

If we want to model arbitrarily long dependencies and hierarchical structures, we have to leave the space of finite-state languages behind and develop formalisms capable of modeling more complex languages. The next section explores the first of such frameworks: context-free languages with the computational models designed to model them.

4.2 Pushdown Language Models

An strong limitation of finite-state language models is that they can definitionally only distinguish a finite set of contexts. However, human language has inherently more structure than what a finite set of contexts can encode. For example, human language contains arbitrarily deep recursive structures which cannot be captured by a finite set of possible histories—we will see an example of this soon in §4.2.1.

To be able to model these structures we are climbing a rung higher on the ladder of the hierarchy of formal languages: we are going to consider **context-free languages**, a larger class of languages than regular languages. Luckily, we will see that a lot of the formal machinery we introduce in this section closely follows analogs from the finite-state section and we invite the reader to pay close attention to the parallels. For example, similarly to how we weighted a string in a regular language by summing over the weights of the paths labeled with that string, we will weight strings in context-free languages by summing over analogous structures.

To be able to recognize context-free languages, will have to extend finite-state automata from §4.1.1 with an additional data structure—the stack. Finite-state automata augmented with a stack are called pushdown automata. We introduce them in §4.2.7. Before giving a formal treatment of pushdown automata, however, we will discuss an arguably more natural formalism for generating the context-free languages—context-free grammars.²¹

In the last part of the section, we will then further extend the regular pushdown automaton with an *additional* stack. Interestingly, this will make it more powerful: as we will see, it will raise its expressive power from context-free languages to all computable languages, as it is Turing complete. While this augmentation will not be immediately useful from a language modeling perspective, we will then later use this machine to prove some theoretical properties of other modern language models we consider later in the course.

4.2.1 Human Language is not Finite-state

As hinted above, human language contains structures that cannot be modeled by finite-state automata. Before we introduce ways of modeling context-free languages, let us, therefore, first motivate the need for a more expressive formalism by more closely considering a specific phenomenon often found in human language: recursive hierarchical structure. We discuss it through an example, based on Jurafsky and Martin (2009).

Example 4.2.1 (Center embeddings). *Consider the sentence:*

“The cat likes to cuddle.”

It simply describes a preference of a cat. However, we can also extend it to give additional information about the cat:

“The cat the dog barked at likes to cuddle.”

This sentence, in turn, can be extended to include additional information about the dog:

“The cat the dog the mouse startled barked at likes to cuddle.”

Of course, we can continue on:

“The cat the dog the mouse the rat frightened startled barked at likes to cuddle.”

and on:

“The cat the dog the mouse the rat the snake scared frightened startled barked at likes to cuddle.”

*In theory, we could continue like this for as long as we wanted—all these sentences are grammatically correct—this is an instance of the so-called **center embeddings**.*

Crucially, such sentences cannot be captured by a regular language, i.e., a language based on an automaton with finitely many states. While we would need formal machinery beyond the scope of this course to formally

²¹You might wonder what non-context-free grammars are: a superclass of context-free grammars is that of context-sensitive grammars, in which a production rule may be surrounded by a left and right context. They are still however a set of restricted cases of general grammars, which are grammars that can emulate Turing machines.

prove this, the intuition is quite simple. By adding more and more “levels” of recursion to the sentences (by introducing more and more animals in the chain), we unboundedly increase the amount of information the model has to “remember” about the initial parts of the sentence while processing it sequentially, to be able to process or generate the matching terms on the other end of the sentence correctly. Because such hierarchies can be arbitrarily deep (and thus the sentences arbitrarily long), there is no bound on the number of states needed to remember them, which means they cannot be captured by a finite-state automaton.

Note that this example also touches upon the distinction of the grammatical competence versus grammatical performance (Chomsky, 1959; Chomsky and Schützenberger, 1963; Chomsky, 1965). The former refers to the purely theoretical properties of human language, for example, the fact that such hierarchical structures can be arbitrarily long and still grammatically correct. Grammatical performance, on the other hand, studies language grounded more in the way people actually use it. For example, nested structures like the one above are never very deep in day-to-day speech—indeed, you probably struggled to understand the last few sentences above. We rarely come across nestings of depth more than three in human language (Miller and Chomsky, 1963; Jin et al., 2018; Karlsson, 2007).

4.2.2 Context-free Grammars

How can we capture recursive structures like those in Example 4.2.1 and the long-term dependencies arising from them? The first formalism modeling such phenomena we will introduce is context-free grammars: a *generative* formalism which can tell us how to generate or “compose” strings in the language it describes. Later in the section (§4.2.7), we will introduce the context-free analog of finite-state automata, which will tell us how to *recognize* whether a string is in a context-free language (rather than generate a string): pushdown automata.

Definition 4.2.1 (Context-free Grammar). A **context-free grammar** (CFG) is a 4-tuple $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P})$ where Σ is an alphabet of terminal symbols, \mathcal{N} is a non-empty set of non-terminal symbols with $\mathcal{N} \cap \Sigma = \emptyset$, $S \in \mathcal{N}$ is the designated start non-terminal symbol and \mathcal{P} is the set of production rules, where each rule $p \in \mathcal{P}$ is of the form $X \rightarrow \alpha$ with $X \in \mathcal{N}$ and $\alpha \in (\mathcal{N} \cup \Sigma)^*$.²²

Example 4.2.2 (A simple context-free grammar). Let $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P})$ be defined as follows:

- $\Sigma = \{a, b\}$
- $\mathcal{N} = \{X\}$
- $S = X$
- $\mathcal{P} = \{X \rightarrow aXb, X \rightarrow \varepsilon\}$

This defines a simple context-free grammar. We will return to it later, when we will formally show that it generates the language $L = \{a^n b^n \mid n \in \mathbb{N}_{\geq 0}\}$.

Rule Applications and Derivations

Context-free grammars allow us to generate strings $y \in \Sigma^*$ by *applying* production rules on its non-terminals. We apply a production rule $X \rightarrow \alpha$ to $X \in \mathcal{N}$ in a rule p by taking X on the right-hand side of p and replacing it with α .²³

Definition 4.2.2 (Rule Application). A production rule $Y \rightarrow \beta, \beta \in (\mathcal{N} \cup \Sigma)^*$, is **applicable** to Y in a rule p , if p takes the form

$$X \rightarrow \alpha Y \gamma, \quad \alpha, \gamma \in (\mathcal{N} \cup \Sigma)^*.$$

The **result of applying** $Y \rightarrow \beta$ to $\alpha Y \gamma$ is $\alpha \beta \gamma$.

²²As is the case for initial states in FSAs, multiple start symbols could be possible. However we consider only one for the sake of simplicity.

²³We say that X is on the right-hand side of a rule p if p takes the form $p = (Y \rightarrow \alpha X \gamma)$, where $\alpha, \gamma \in (\mathcal{N} \cup \Sigma)^*$. We will sometimes refer to X as the *head* of the production rule $X \rightarrow \alpha$, and the right-hand side α as the *body* of the production rule.

Starting with S , we apply $S \rightarrow \alpha$ to S for some $(S \rightarrow \alpha) \in \mathcal{P}$, then take a non-terminal in α and apply a new production rule.²⁴ To generate a string we follow this procedure until all non-terminal symbols have been transformed into terminal symbols. The resulting string, i.e., the **yield**, will be the string taken by concatenating all terminal symbols read from left to right. More formally, a derivation can be defined as follows.

Definition 4.2.3 (Derivation). *A **derivation** in a grammar \mathcal{G} is a sequence $\alpha_1, \dots, \alpha_M$, where $\alpha_1 \in \mathcal{N}$, $\alpha_2, \dots, \alpha_{M-1} \in (\mathcal{N} \cup \Sigma)^*$ and $\alpha_M \in \Sigma^*$, in which each α_{m+1} is formed by applying a production rule in \mathcal{P} to α_m .*

We say that $\alpha \in (\mathcal{N} \cup \Sigma)^*$ is derived from $X \in \mathcal{N}$ if we can apply a finite sequence of production rules to generate α starting from X . We will denote this as $X \xRightarrow{*}_{\mathcal{G}} \alpha$. See the following formal definition.

Definition 4.2.4 (Derives). *Let $\mathcal{G} \stackrel{\text{def}}{=} (\Sigma, \mathcal{N}, S, \mathcal{P})$ be a CFG. We say that X **derives** β under the grammar \mathcal{G} , denoted as $X \Rightarrow_{\mathcal{G}} \beta$ if $\exists p \in \mathcal{P}$ such that $p = (X \rightarrow \alpha\beta\gamma)$, $\alpha, \gamma \in (\mathcal{N} \cup \Sigma)^*$ and $\beta \in (\mathcal{N} \cup \Sigma)^* \setminus \{\varepsilon\}$. The special case $X \Rightarrow_{\mathcal{G}} \varepsilon$ holds iff $X \rightarrow \varepsilon \in \mathcal{P}$. We denote the reflexive transitive closure of the $\Rightarrow_{\mathcal{G}}$ relation as $\xRightarrow{*}_{\mathcal{G}}$. We say that β is **derived from** X if $X \xRightarrow{*}_{\mathcal{G}} \beta$.*

The (context-free) language of a CFG \mathcal{G} is defined as all the strings $y \in \Sigma^*$ that can be derived from the start symbol S of \mathcal{G} , or alternatively, the set of all yields possible from derivations in \mathcal{G} that start with S . We will denote the language generated by \mathcal{G} as $L(\mathcal{G})$.

Definition 4.2.5 (Language of a Grammar). *The **language** of a context-free grammar \mathcal{G} is*

$$L(\mathcal{G}) = \{y \in \Sigma^* \mid S \xRightarrow{*}_{\mathcal{G}} y\} \quad (4.84)$$

Parse Trees and Derivation Sets

A natural representation of a derivation in a context-free grammar is a **derivation tree** d (also known as a parse tree). A derivation tree represents the sequence of applied rules in a derivation with a directed tree. The tree's internal nodes correspond to the non-terminals in the derivation, and each of their children corresponds to a symbol (from $\Sigma \cup \mathcal{N}$) on the right side of the applied production in the derivation. The leaves, representing terminal symbols, “spell out” the derived string—the tree's yield. More formally, for each production rule $X \rightarrow \alpha$, the node corresponding to the specific instance of the non-terminal X in the derivation is connected to the nodes corresponding to $Y \in \alpha$ where $Y \in \Sigma \cup \mathcal{N}$.

We will mostly be interested in representing derivations starting with S —the root node of a tree representing any such derivation will correspond to S . We will denote the string generated by a tree d —its yield—by $s(d)$. See Fig. 4.5 for examples of parse trees for the grammar from Example 4.2.2.

Importantly, a grammar may in fact admit *multiple* derivations and hence multiple derivation trees for any given string.

Example 4.2.3 (Multiple derivation strings). *It is relatively easy to see that in the grammar \mathcal{G} , each string $a^n b^n$ is only generated by a single derivation tree—each new pair of symbols a and b can only be added by applying the rule $X \rightarrow aXb$ and the string $a^n b^n$ can only be generated by the application of the rule $X \rightarrow aXb$ n times and the rule $X \rightarrow \varepsilon$ once in this order.*

However, we can modify \mathcal{G} by adding, for instance, a non-terminal Y and rules $X \rightarrow Y, Y \rightarrow \varepsilon$. The empty string ε may then be derived either by $(X \rightarrow \varepsilon)$, or $(X \rightarrow Y), (Y \rightarrow \varepsilon)$, corresponding to two separate derivation trees, as shown in Fig. 4.6. The set of these two trees comprises what we call the derivation set of ε .

We denote a derivation set of a string y , generated by the grammar \mathcal{G} , as $\mathcal{D}_{\mathcal{G}}(y)$.

Definition 4.2.6 (String derivation set). *Let $y \in \Sigma^*$. Its **derivation set**, denoted by $\mathcal{D}_{\mathcal{G}}(y)$ is defined as*

$$\mathcal{D}_{\mathcal{G}}(y) \stackrel{\text{def}}{=} \{d \mid s(d) = y\}. \quad (4.85)$$

²⁴We will write $X \in \alpha$, which formally means a substring of α with length 1. Unless otherwise stated, X can be either a non-terminal or a terminal.

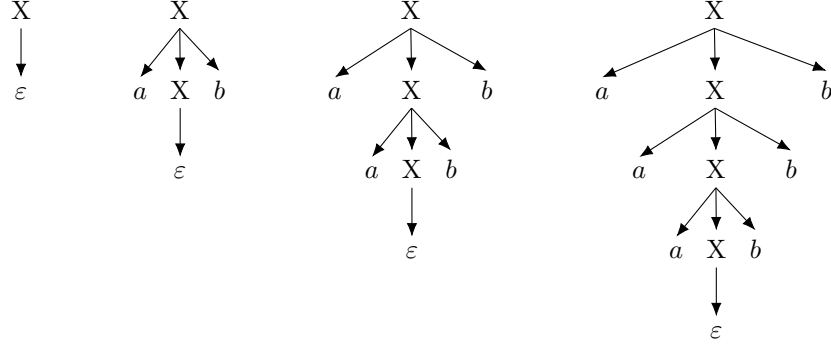


Figure 4.5: A sequence of derivation trees for the strings in $\{a^n b^n \mid n = 0, 1, 2, 3\}$ in the grammar from Example 4.2.2.

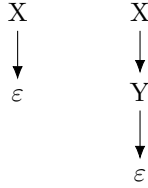


Figure 4.6: Two parse trees in the modified grammar \mathcal{G} yielding ε .

We say that a grammar is **unambiguous** if, for every string that can be generated by the grammar, there is only one associated derivation tree.

Definition 4.2.7 (Unambiguity). *A grammar \mathcal{G} is **unambiguous** if for all $y \in L(\mathcal{G})$, $|\mathcal{D}_{\mathcal{G}}(y)| = 1$.*

The converse holds for **ambiguous** grammars.

Definition 4.2.8 (Ambiguity). *A grammar \mathcal{G} is **ambiguous** if $\exists y \in L(\mathcal{G})$ such that $|\mathcal{D}_{\mathcal{G}}(y)| > 1$.*

The set of all derivation trees in a grammar is its derivation set.

Definition 4.2.9 (Grammar derivation set). *The **derivation set of a grammar**, $\mathcal{D}_{\mathcal{G}}$, is the set of all derivations possible under the grammar. More formally, it can be defined as the union over the derivation set for the strings in its language,*

$$\mathcal{D}_{\mathcal{G}} \stackrel{\text{def}}{=} \bigcup_{y' \in L(\mathcal{G})} \mathcal{D}_{\mathcal{G}}(y') \quad (4.86)$$

Definition 4.2.10 (Non-terminal derivation set). *The **derivation set of a non-terminal** $Y \in \mathcal{N}$ in \mathcal{G} , denoted $\mathcal{D}_{\mathcal{G}}(Y)$, is defined as the set of derivation subtrees with root node Y .*

Note that $\mathcal{D}_{\mathcal{G}}$ could be defined as $\mathcal{D}_{\mathcal{G}}(S)$. For a terminal symbol $a \in \Sigma$, we trivially define the derivation set $\mathcal{D}_{\mathcal{G}}(a)$ to be empty.²⁵

In cases where it is irrelevant to consider the order of the production rules in a derivation tree, we will write $(X \rightarrow \alpha) \in \mathbf{d}$ to refer to specific production rules in the tree—viewing trees as multisets (or bags) over the production rules they include.

Example 4.2.4 (Nominal Phrases). *CFGs are often used to model natural languages. Terminals would then correspond to words in the natural language, strings would be text sequences and non-terminals would be abstractions over words. As an example, consider a grammar \mathcal{G} that can generate a couple of nominal phrases. We let $\mathcal{N} = \{\text{Adj}, \text{Det}, \text{N}, \text{Nominal}, \text{NP}\}$, $\Sigma = \{a, \text{big}, \text{female}, \text{giraffe}, \text{male}, \text{tall}, \text{the}\}$, $S = \text{Nominal}$ and define the following production rules:*

²⁵Empty derivation sets for terminal symbols is defined solely for ease of notation later.

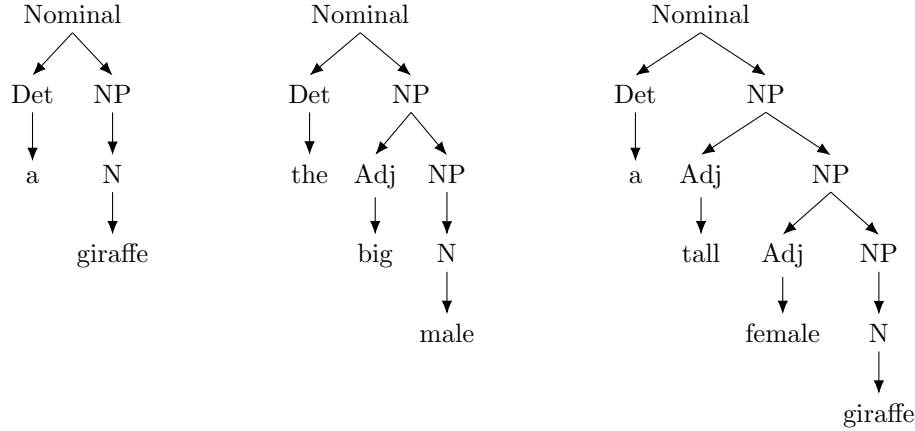


Figure 4.7: Derivation trees for natural language nominal phrases.

$\text{Nominal} \rightarrow \text{Det NP}$
 $\text{NP} \rightarrow \text{N} \mid \text{Adj NP}$
 $\text{Det} \rightarrow a \mid the$
 $\text{N} \rightarrow female \mid giraffe \mid male$
 $\text{Adj} \rightarrow big \mid female \mid male \mid tall$

See Fig. 4.7 for a few examples of derivation trees in this grammar.

Example 4.2.5 (The generalized Dyck languages $D(k)$). A very widely studied family of context-free languages are the Dyck- k languages, $D(k)$, the languages of well-nested brackets of k types. They are, in some ways, archetypal context-free languages (Chomsky and Schützenberger, 1963). Formally, we can define them as follows.

Definition 4.2.11 ($D(k)$ languages). Let $k \in \mathbb{N}$. The $D(k)$ language is the language of the following context-free grammar $\mathcal{G} \stackrel{\text{def}}{=} (\Sigma, \mathcal{N}, S, \mathcal{P})$

- $\Sigma \stackrel{\text{def}}{=} \{\langle_n \mid n = 1, \dots, k\} \cup \{\rangle_n \mid n = 1, \dots, k\}$
- $\mathcal{N} \stackrel{\text{def}}{=} \{S\}$
- $S \stackrel{\text{def}}{=} S$
- $\mathcal{P} \stackrel{\text{def}}{=} \{S \rightarrow \varepsilon, S \rightarrow SS\} \cup \{S \rightarrow \langle_n S \rangle_n \mid n = 1, \dots, k\}$

Examples of strings in the language $D(3)$ would be $\langle_3 \rangle_3 \langle_2 \rangle_2 \langle_1 \rangle_1$, $\langle_3 \rangle_3 \langle_1 \rangle_1 \langle_2 \rangle_2 \langle_2 \rangle_2 \langle_1 \rangle_1$, and $\langle_1 \rangle_1 \langle_2 \rangle_2 \langle_2 \rangle_2 \langle_2 \rangle_2 \langle_3 \rangle_3 \langle_1 \rangle_1 \langle_3 \rangle_3$. The string $\langle_2 \rangle_2 \langle_1 \rangle_1 \rangle_2$ is not in the language $D(3)$.

To give you a taste of what formally working with context-free grammars might look like, we now formally show that the grammar from Example 4.2.2 really generates the language $L = \{a^n b^n \mid n \in \mathbb{N}_{\geq 0}\}$, as we claimed.

Example 4.2.6 (Recognizing $a^n b^n$). The language $L = \{a^n b^n \mid n \in \mathbb{N}\}$ is not regular.²⁶ However, we can show that it is context-free and recognized exactly by the simple grammar from Example 4.2.2. We restate it here for convenience: $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P})$ with $\mathcal{N} = \{X\}$, $\Sigma = \{a, b\}$, $S = X$, $\mathcal{P} = \{X \rightarrow aXb, X \rightarrow \varepsilon\}$.

Lemma 4.2.1. Given the grammar \mathcal{G} defined above, we have $L(\mathcal{G}) = \{a^n b^n \mid n \in \mathbb{N}\}$.

Proof. We will show that $L = L(\mathcal{G})$ in two steps: (i) showing that $L \subseteq L(\mathcal{G})$ and (ii) showing that $L(\mathcal{G}) \subseteq L$. Define $\mathbf{y}_n = a^n b^n$.

(i) We first need to show that each $\mathbf{y} \in L$ can be generated by \mathcal{G} , which we will do by induction.

²⁶Again, while the intuition behind this is similar to our reasoning from Example 4.2.1, this would have to be proven using the so-called pumping lemma for regular languages.

Base case ($n = 0$) We have that $y_0 = \varepsilon$, which is generated by $d = (X \rightarrow \varepsilon)$.

Inductive step ($n > 1$) We have that y_n is generated by

$$d = \underbrace{(X \rightarrow aXb) \cdots (X \rightarrow aXb)}_{n \text{ times}} (X \rightarrow \varepsilon).$$

It is then easy to see that y_{n+1} is generated by the derivation we get by replacing the last rule $(X \rightarrow \varepsilon)$ with $(X \rightarrow aXb)(X \rightarrow \varepsilon)$ —they are exactly the trees illustrated in Fig. 4.5.

(ii) Next, we show that for each $d \in \mathcal{D}_G$, we have that $y(d) \in L$.

Base case ($d = (X \rightarrow \varepsilon)$) It is trivial to see that the derivation $d = (X \rightarrow \varepsilon)$ yields $y(d) = \varepsilon$.

Inductive step Now observe that \mathcal{P} only contains two production rules and one non-terminal. Starting with X , we can either apply $X \rightarrow aXb$ to get one new non-terminal X , or apply $X \rightarrow \varepsilon$ to terminate the process. Hence, if we fix the length of the sequence of production rules, there is no ambiguity in which string will be generated. Thus, by induction, we conclude that if we have a derivation tree given by $\underbrace{(X \rightarrow aXb), \dots, (X \rightarrow aXb)}_{n \text{ times}}, (X \rightarrow \varepsilon)$ generating $a^n b^n$, the derivation tree given by $\underbrace{(X \rightarrow aXb), \dots, (X \rightarrow aXb)}_{n+1 \text{ times}}, (X \rightarrow \varepsilon)$ will generate $a^{n+1} b^{n+1}$. ■

Reachable Non-terminals and Pruning

Similarly to how some states in a WFSA can be useless in the sense that they are not accessible from an initial state or might not lead to a final state, so too can non-terminals in a CFG be useless by not being reachable from the start symbol or might not lead to any string of terminals. In the context of CFGs, we typically use a different terminology: “reachable” instead of “accessible” and “generating” instead of “co-accessible”.

Definition 4.2.12 (Accessibility for CFGs). A symbol $X \in \mathcal{N} \cup \Sigma$ is **reachable** (or **accessible**) if $\exists \alpha, \alpha' \in (\mathcal{N} \cup \Sigma)^*$ such that $S \xRightarrow{*} \alpha X \alpha'$.

Definition 4.2.13 (Co-accessibility for CFGs). A non-terminal Y is **generating** (or **co-accessible**) if $\exists y \in \Sigma^*$ such that $Y \xRightarrow{*} y$.

In words, reachable symbols are those that can be derived from the start symbol, whereas generating non-terminals are those from which at least one string (including the empty string) can be derived. Note that we define reachable for both non-terminals and terminals while generating is only defined for non-terminals.

This allows us to define a pruned context-free grammar, which is the CFG version of a trimmed WFSA.

Definition 4.2.14 (Pruned CFG). A CFG is **pruned** (or **trimmed**) if it has no useless non-terminals, i.e. all non-terminals are both reachable and generating. **Pruning** (or **trimming**) refers to the removal of useless non-terminals.

4.2.3 Weighted Context-free Grammars

As we did with finite-state automata, we will augment the classic, unweighted context-free grammars with real-valued weights. We do that by associating with each rule $X \rightarrow \alpha$ a weight $\mathcal{W}(X \rightarrow \alpha) \in \mathbb{R}$.

Definition 4.2.15 (Weighted Context-free Grammar). A **real-weighted context-free grammar** is a 5-tuple $(\Sigma, \mathcal{N}, S, \mathcal{P}, \mathcal{W})$ where Σ is an alphabet of terminal symbols, \mathcal{N} is a non-empty set of non-terminal symbols with $\mathcal{N} \cap \Sigma = \emptyset$, $S \in \mathcal{N}$ is the designated start non-terminal symbol, \mathcal{P} is the set of production rules, and \mathcal{W} a function $\mathcal{W}: \mathcal{P} \rightarrow \mathbb{R}$, assigning each production rule a real-valued weight.

For notational brevity, we will denote rules $p \in \mathcal{P}$ as $p = X \xrightarrow{w} \alpha$ for $X \in \mathcal{N}$, $\alpha \in (\mathcal{N} \cup \Sigma)^*$ and $w = \mathcal{W}(X \rightarrow \alpha) \in \mathbb{R}$.

Example 4.2.7 (A simple weighted context-free grammar). Consider the grammar $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P}, \mathcal{W})$ defined as follows:

- $\Sigma = \{a, b\}$
- $\mathcal{N} = \{X\}$
- $S = X$
- $\mathcal{P} = \{X \rightarrow aXb, X \rightarrow \varepsilon\}$
- $\mathcal{W} = \{X \rightarrow aXb \mapsto \frac{1}{2}, X \rightarrow \varepsilon \mapsto \frac{1}{2}\}$

This defines a simple weighting of the CFG from Example 4.2.2.

Weights assigned to productions by WFCGs can be arbitrary real numbers. Analogous to probabilistic WFSA (Definition 4.1.17) describing locally normalized finite-state language models, we also define probabilistic WCFGs, where the weights of *applicable production rules* to any non-terminal form a probability distribution.

Definition 4.2.16. A weighted context-free grammar $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P}, \mathcal{W})$ is **probabilistic** if the weights of the productions of every non-terminal are non-negative and sum to 1, i.e., for all $X \in \mathcal{N}$, it holds that

$$\forall X \rightarrow \alpha \in \mathcal{P}, \mathcal{W}(X \rightarrow \alpha) \geq 0 \quad (4.87)$$

and

$$\sum_{X \rightarrow \alpha \in \mathcal{P}} \mathcal{W}(X \rightarrow \alpha) = 1 \quad (4.88)$$

Intuitively, this means that all the production weights are non-negative and that, for any left side of a production rule X , the weights over all production rules $X \rightarrow \alpha$ sum to 1. The grammar from Example 4.2.7 is, therefore, also probabilistic.

Again analogously to the WFSA case, we say that a string \mathbf{y} is in the language of WCFG \mathcal{G} if there exists a derivation tree \mathbf{d} in \mathcal{G} containing only non-zero weights with yield $\mathbf{s}(\mathbf{d}) = \mathbf{y}$.

Tree Weights, String Weights, and Allsums

In the case of regular languages, we discussed how individual strings are “produced” by paths in the automaton (in the sense that each path *yields* a string). As Example 4.2.4 showed, the structures that “produce” or yield strings in a context-free grammar are *trees*—those, therefore, play an analogous role in context-free grammars to paths in finite-state automata.

Just like we asked ourselves how to combine individual transition weights in a WFSA into weights of entire paths and later how to combine those into weights of strings, we now consider the questions of how to combine the weights of individual production rules into the weight of entire trees and later also individual strings. We start by giving a definition of the weight of a tree as the product over the weights of all the rules in the tree, i.e., as a *multiplicatively decomposable* function over the weights of its rules. As you can probably foresee, we will then define the weight of a string as the *sum* over all the trees which yield that string.

Definition 4.2.17. The **weight of a derivation tree** $\mathbf{d} \in \mathcal{D}_{\mathcal{G}}$ defined by a WCFG \mathcal{G} is

$$\mathbf{w}(\mathbf{d}) = \prod_{(X \rightarrow \alpha) \in \mathbf{d}} \mathcal{W}(X \rightarrow \alpha). \quad (4.89)$$

The stringsum or the string acceptance weight of a particular string under a grammar is then defined as follows:

Definition 4.2.18. The **stringsum** $\mathcal{G}(\mathbf{y})$ of a string \mathbf{y} generated by a WCFG \mathcal{G} is defined by

$$\mathcal{G}(\mathbf{y}) = \sum_{\mathbf{d} \in \mathcal{D}_{\mathcal{G}}(\mathbf{y})} \mathbf{w}(\mathbf{d}) \quad (4.90)$$

$$= \sum_{\mathbf{d} \in \mathcal{D}_{\mathcal{G}}(\mathbf{y})} \prod_{(X \rightarrow \alpha) \in \mathbf{d}} \mathcal{W}(X \rightarrow \alpha) \quad (4.91)$$

Lastly, analogously to the allsum in WFSAs, an **allsum** is the sum of the weights of all the trees in a WCFG. We first define the allsum for symbols (non-terminals and terminals).

Definition 4.2.19 (Allsum). *The **allsum** for a non-terminal Y in a grammar \mathcal{G} is defined by*

$$Z(\mathcal{G}, Y) = \sum_{\mathbf{d} \in \mathcal{D}_{\mathcal{G}}(Y)} w(\mathbf{d}) \quad (4.92)$$

$$= \sum_{\mathbf{d} \in \mathcal{D}_{\mathcal{G}}(Y)} \prod_{(X \rightarrow \alpha) \in \mathbf{d}} W(X \rightarrow \alpha) \quad (4.93)$$

The allsum for a terminal $a \in \Sigma \cup \{\varepsilon\}$ is defined to be

$$Z(a) \stackrel{\text{def}}{=} \mathbf{1}. \quad (4.94)$$

The allsum for a grammar is then simply the allsum for its start symbol.

Definition 4.2.20 (Allsum WCFG). *The **allsum** of a weighted context-free grammar $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P}, W)$ is*

$$Z(\mathcal{G}) = Z(\mathcal{G}, S) \quad (4.95)$$

$$= \sum_{\mathbf{d} \in \mathcal{D}_{\mathcal{G}}(S)} w(\mathbf{d}) \quad (4.96)$$

$$= \sum_{\mathbf{d} \in \mathcal{D}_{\mathcal{G}}(S)} \prod_{(X \rightarrow \alpha) \in \mathbf{d}} W(X \rightarrow \alpha) \quad (4.97)$$

When the grammar \mathcal{G} we refer to is clear from context, we will drop the subscript and write e.g. $Z(S)$.

Although we can in some cases compute the allsum of a WCFG in closed form, as we will see in the example below, we generally require some efficient algorithm to be able to do so.

Example 4.2.8 (Geometric Series as an Allsum). *Consider the WCFG $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P}, W)$, given by $\mathcal{N} = \{X\}$, $\Sigma = \{a\}$, $S = X$, and the rules:*

$$\begin{aligned} X &\xrightarrow{1/3} aX \\ X &\xrightarrow{1} \varepsilon \end{aligned}$$

The language generated by \mathcal{G} is $L(\mathcal{G}) = \{a^n \mid n \geq 0\}$. Further note that this grammar is unambiguous – each string $\mathbf{y} = a^m$, for some $m \geq 0$, is associated with the derivation tree given by $\underbrace{(X \xrightarrow{1/3} aX), \dots, (X \xrightarrow{1/3} aX)}_{m \text{ times}}, (X \xrightarrow{1} \varepsilon)$.

Due to the multiplicative decomposition over the weights of the rules, the weight associated with each derivation tree \mathbf{d} will hence be

$$w(\mathbf{d}) = \left(\frac{1}{3}\right)^m \times 1 = \left(\frac{1}{3}\right)^m$$

Accordingly, we can compute the allsum of \mathcal{G} using the closed-form expression for geometric series:

$$Z(\mathcal{G}) = \sum_{m=0}^{\infty} \left(\frac{1}{3}\right)^m = \frac{1}{1 - 1/3} = \frac{3}{2}$$

Just like we defined normalizable WFSAs, we also define normalizable WCFGs in terms of their allsum.

Definition 4.2.21 (Normalizable Weighted Context-free Grammar). *A weighted context-free grammar \mathcal{G} is **normalizable** if $Z(\mathcal{G})$ is finite, i.e., $Z(\mathcal{G}) < \infty$.*

4.2.4 Context-free Language Models

This brings us to the definition of context-free language models.

Definition 4.2.22. *A language model p_{LM} is **context-free** if its weighted language equals the language of some weighted context-free grammar, i.e., if there exists a weighted context-free grammar \mathcal{G} such that $L(\mathcal{G}) = L(p_{LM})$.*

Going the other way—defining string probabilities given a weighted context-free grammar—there are again two established ways of defining the probability of a string in its language.

String Probabilities in a Probabilistic Context-free Grammar

In a probabilistic CFG (cf. Definition 4.2.16), any production from a non-terminal $X \in \mathcal{N}$ is associated with a probability. As the probabilities of continuing a derivation (and, therefore, a derivation tree) depend solely on the individual terminals (this is the core of *context-free* grammars!), it is intuitive to see those probabilities as conditional probabilities of the new symbols given the output generated so far. One can, therefore, define the probability of a path as the product of these individual “conditional” probabilities.

Definition 4.2.23 (Tree probability in a PCFG). *We call the weight of a tree $d \in \mathcal{D}_{\mathcal{G}}$ in a probabilistic CFG the **probability** of the tree d .*

This alone is not enough to define the probability of any particular string $y \in \Sigma^*$ since there might be multiple derivations of y . Naturally, we define the probability of y as the sum of the individual trees that generate it:

Definition 4.2.24 (String probability in a PCFG). *We call the stringsum of a string $y \in \Sigma^*$ in a probabilistic CFG \mathcal{G} the **probability** of the string y :*

$$p_{\mathcal{G}}(y) \stackrel{\text{def}}{=} \mathcal{G}(y). \quad (4.98)$$

These definitions and their affordances mirror the ones in probabilistic finite-state automata (cf. §4.1.2): they again do not require any *normalization* and are therefore attractive as the summation over all possible strings is avoided. Again, the question of *tightness* of such models comes up: we explore it question in §4.2.5.

String Probabilities in a General Weighted Context-free Grammar

To define string probabilities in a general weighted CFG, we use the introduced notions of the stringsum and the allsum—we *normalize* the stringsum to define the globally normalized probability of a string y as the *proportion* of the total weight assigned to all strings that is assigned to y .

Definition 4.2.25 (String probability in a WCFG). *Let $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P}, \mathcal{W})$ be a normalizable WCFG with non-negative weights. We define the probability of a string $y \in \Sigma^*$ under \mathcal{G} as*

$$p_{\mathcal{G}}(y) \stackrel{\text{def}}{=} \frac{\mathcal{G}(y)}{Z(\mathcal{G})}. \quad (4.99)$$

Language Models Induced by a Weighted Context-free Grammar

With the notions of string probabilities in both probabilistic and general weighted CFGs, we can now define the language model induced by \mathcal{G} as follows.

Definition 4.2.26 (A language model induced by a WCFG). *Let $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P}, \mathcal{W})$ be a WCFG. We define the **language model induced by \mathcal{G}** as the following probability distribution over Σ^**

$$p_{LM\mathcal{G}}(y) \stackrel{\text{def}}{=} p_{\mathcal{G}}(y). \quad (4.100)$$

Again, it is easy to see that while global normalization requires the computation of the allsum, language models induced by weighted FSAs through Eq. (4.99) are *globally normalized* and thus always tight. The tightness of *probabilistic* WCFGs is discussed next, after which we investigate the relationship between globally- and locally-normalized context-free grammars.

4.2.5 Tightness of Context-free Language Models

Again, an advantage of globally normalized context-free language models (grammars) is that they are always tight, as the derivation trees are explicitly normalized with the global normalization constant such that they sum to 1 over the set of possible sentences.

In this section, we, therefore, consider the tightness of probabilistic context-free grammars. We follow the exposition from [Booth and Thompson \(1973\)](#). The proof requires the use of multiple new concepts, which we first introduce below.

Definition 4.2.27 (Generation level). *We define the **level of a generation sequence** inductively as follows. The zeroth level γ_0 of a generation sequence is defined as S . Then, for any $n > 0$, γ_n corresponds to the string is obtained by applying the applicable productions onto all nonterminals of γ_{n-1} .*

Example 4.2.9. Let $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P})$ with $\Sigma = \{a, b\}$, $\mathcal{N} = \{S, X, Y\}$, and $\mathcal{P} = \{S \rightarrow aXY, X \rightarrow YX, X \rightarrow bYY, Y \rightarrow aY, Y \rightarrow a\}$. Then the generation sequence of the string *aabaaaa* would be

$$\begin{aligned} \gamma_0 &= S && \text{(definition)} \\ \gamma_1 &= aXY && \text{(applying } S \rightarrow aXY) \\ \gamma_2 &= aYXaaY && \text{(applying } X \rightarrow YX, Y \rightarrow aY) \\ \gamma_3 &= aabYYaaaaY && \text{(applying } Y \rightarrow a, X \rightarrow bYY, X \rightarrow aY) \\ \gamma_3 &= aabaaaaaa && \text{(applying } Y \rightarrow a, Y \rightarrow a, Y \rightarrow a) \end{aligned}$$

We will also rely heavily on generating functions. A generating function is simply a way of *representing* an infinite sequence by encoding its elements as the coefficients of a *formal power series*. Unlike ordinary series such as the geometric power series from [Example 4.2.8](#), a formal power series does not need to converge: in fact, at its core a generating function is not actually regarded as a *function*—its “variables” are indeterminate and they simply serve as “hooks” for the numbers in the sequence.

Definition 4.2.28 (Production generating function). Let $\mathcal{G} \stackrel{\text{def}}{=} (\Sigma, \mathcal{N}, S, \mathcal{P}, \mathcal{W})$ be a PCFG and $N \stackrel{\text{def}}{=} |\mathcal{N}|$. For each $X_n \in \mathcal{N}$, define its **production generating function** as

$$g(s_1, \dots, s_N) \stackrel{\text{def}}{=} \sum_{X_n \rightarrow \alpha} \mathcal{W}(X_n \rightarrow \alpha) s_1^{r_1(\alpha)} s_2^{r_2(\alpha)} \dots s_N^{r_N(\alpha)}, \quad (4.101)$$

where $r_m(\alpha)$ denotes the number of times the nonterminal $X_m \in \mathcal{N}$ appears in $\alpha \in (\Sigma \cup \mathcal{N})^*$.

Example 4.2.10. Let $\mathcal{G} = (\Sigma, \mathcal{N}, S, \mathcal{P})$ with $\Sigma = \{a, b\}$, $\mathcal{N} = \{S, X\}$, and $\mathcal{P} = \{S \rightarrow aSX, S \rightarrow b, X \rightarrow aXX, X \rightarrow aa\}$. Then

$$\begin{aligned} g_1(s_1, s_2) &= \mathcal{W}(S \rightarrow aSX) s_1 s_2 + \mathcal{W}(S \rightarrow b) \\ g_2(s_1, s_2) &= \mathcal{W}(X \rightarrow aXX) s_2^2 + \mathcal{W}(X \rightarrow aa) \end{aligned}$$

Definition 4.2.29 (Generating function). The **generating function** of the l^{th} level is defined as

$$G_0(s_1, \dots, s_N) \stackrel{\text{def}}{=} s_1 \quad (4.102)$$

$$G_1(s_1, \dots, s_N) \stackrel{\text{def}}{=} g_1(s_1, \dots, s_N) \quad (4.103)$$

$$G_l(s_1, \dots, s_N) \stackrel{\text{def}}{=} G_{l-1}(g_1(s_1, \dots, s_N), \dots, g_N(s_1, \dots, s_N)), \quad (4.104)$$

that is, the l^{th} -level generating function is defined as the $l - 1^{\text{st}}$ -level generating function applied to production generating functions as arguments.

Example 4.2.11. For the grammar from Example 4.2.10, we have

$$\begin{aligned}
G_0(s_1, s_2) &= s_1 \\
G_1(s_1, s_2) &= g(s_1, s_2) = \mathcal{W}(S \rightarrow aSX) s_1 s_2 + \mathcal{W}(S \rightarrow b) \\
G_2(s_1, s_2) &= \mathcal{W}(S \rightarrow aSX) [g_1(s_1, s_2)] [g_2(s_1, s_2)] + \mathcal{W}(S \rightarrow b) \\
&= \mathcal{W}(S \rightarrow aSX)^2 \mathcal{W}(X \rightarrow aXX) s_1 s_2^3 \\
&\quad + \mathcal{W}(S \rightarrow aSX)^2 \mathcal{W}(X \rightarrow aa) s_1 s_2 \\
&\quad + \mathcal{W}(S \rightarrow aSX) \mathcal{W}(S \rightarrow b) \mathcal{W}(X \rightarrow aXX) s_2^2 \\
&\quad + \mathcal{W}(S \rightarrow aSX) \mathcal{W}(S \rightarrow b) \mathcal{W}(X \rightarrow aa) \\
&\quad + \mathcal{W}(S \rightarrow b)
\end{aligned}$$

We can see that a generating function $G_l(s_1, \dots, s_N)$ can be expressed as

$$G_l(s_1, \dots, s_N) = D_l(s_1, \dots, s_N) + C_l \quad (4.105)$$

where the polynomial $D_l(s_1, \dots, s_N)$ does not contain any constant terms. It is easy to see that the constant C_l then corresponds to the probability of all strings that can be derived in l levels or fewer. This brings us to the following simple lemma.

Lemma 4.2.2. A PCFG is tight if and only if

$$\lim_{l \rightarrow \infty} C_l = 1. \quad (4.106)$$

Proof. Suppose that $\lim_{l \rightarrow \infty} C_l < 1$. This means that the generation process can enter a generation sequence that has a non-zero probability of not terminating—this corresponds exactly to it not being tight.

On the other hand, $\lim_{l \rightarrow \infty} C_l = 1$ implies that no such sequence exists, since the limit represents the probability of all strings that can be generated by derivations of a finite number of production rules. ■

The rest of the section considers necessary and sufficient conditions for Eq. (4.106) to hold. For this, we first define the first-moment matrix of a PCFG.

Definition 4.2.30. Let $\mathcal{G} \stackrel{\text{def}}{=} (\Sigma, \mathcal{N}, S, \mathcal{P}, \mathcal{W})$ be a PCFG. We define its **first-moment matrix** (its mean matrix) $\mathbf{E} \in \mathbb{R}^{N \times N}$ as

$$E_{nm} \stackrel{\text{def}}{=} \left. \frac{\partial g_n(s_1, \dots, s_N)}{\partial s_m} \right|_{s_1, \dots, s_N=1}. \quad (4.107)$$

Note that E_{nm} represents the *expected number of occurrences* of the non-terminal X_m in the set of sequences α with $X_n \Rightarrow_{\mathcal{G}} \alpha$, i.e., the set of sequences X_n can be rewritten into:

$$E_{nm} = \sum_{X_n \rightarrow \alpha} \mathcal{W}(X_n \rightarrow \alpha) r_m(\alpha). \quad (4.108)$$

The informal intuition behind this is the following: each of the terms $\mathcal{W}(X_n \rightarrow \alpha) s_1^{r_1(\alpha)} s_2^{r_2(\alpha)} \dots s_N^{r_N(\alpha)}$ in g_n contains the information about *how many times* any non-terminal X_m appears in the production rule $X_n \rightarrow \alpha$ as well as what the probability of “using” or applying that production rule to X_n is. Differentiating $\mathcal{W}(X_n \rightarrow \alpha) s_1^{r_1(\alpha)} s_2^{r_2(\alpha)} \dots s_N^{r_N(\alpha)}$ w.r.t. s_m then “moves” the coefficient r_m corresponding to the number of occurrences of X_m in $X_n \rightarrow \alpha$ in front of the term $\mathcal{W}(X_n \rightarrow \alpha) s_1^{r_1(\alpha)} s_2^{r_2(\alpha)} \dots s_N^{r_N(\alpha)}$ in g_n , effectively multiplying the probability of the occurrence of the rule with the number of terms X_m in the rule—this is exactly the expected number of occurrences of X_m for this particular rule, averaging over all possible rules that could be applied. Summing over all applicable production rules for X_n (which form a probability distribution) gives us the total expected number of occurrences of X_m . This brings us to the core theorem of this section characterizing the tightness of PCFGs.

Theorem 4.2.1 (Tightness of PCFGs). A PCFG is tight if $|\lambda_{\max}| < 1$ and is non-tight if $|\lambda_{\max}| > 1$, where λ_{\max} is the eigenvalue of \mathbf{E} with the largest absolute value.

Proof. The coefficient of the term $s_1^{r_1} s_2^{r_2} \cdots s_N^{r_N}$ in the generating function $G_l(s_1, \dots, s_N)$ corresponds to the probability that there will be r_1 non-terminal symbols X_1, \dots, r_N non-terminal symbols X_N in the l^{th} level of the generation sequence. In particular, if the grammar is tight, this means that

$$\lim_{l \rightarrow \infty} G_l(s_1, \dots, s_N) = \lim_{l \rightarrow \infty} [D_l(s_1, \dots, s_N) + C_l] = 1. \quad (4.109)$$

This, however, is only true if

$$\lim_{l \rightarrow \infty} D_l(s_1, \dots, s_N) = 0 \quad (4.110)$$

and this, in turn, can only be true if $\lim_{l \rightarrow \infty} r_n = 0$ for all $n = 1, \dots, N$. The expected value of r_n at level l is

$$\bar{r}_{l,n} = \left. \frac{\partial G_l(s_1, \dots, s_N)}{\partial s_n} \right|_{s_1, \dots, s_N=1}. \quad (4.111)$$

Reasoning about this is similar to the intuition behind the first-moment matrix, with the difference that we are now considering the number of occurrences after a sequence of l applications. Denoting

$$\bar{\mathbf{r}}_l \stackrel{\text{def}}{=} [\bar{r}_{l,1}, \dots, \bar{r}_{l,N}] \quad (4.112)$$

we have

$$\bar{\mathbf{r}}_l = \left[\sum_{j=1}^N \frac{\partial G_{l-1}(g_1(s_1, \dots, s_N), \dots, g_N(s_1, \dots, s_N))}{\partial g_j} \right]_{s_1, \dots, s_N=1} \quad (4.113)$$

$$\cdot \left. \frac{\partial g_j(s_1, \dots, s_N)}{\partial s_n} \right|_{s_1, \dots, s_N=1} \Big|_{n=1, \dots, N} \quad (4.114)$$

$$= \bar{\mathbf{r}}_{l-1} \mathbf{E}. \quad (4.115)$$

Applying this relationship repeatedly, we get

$$\bar{\mathbf{r}}_l = \bar{\mathbf{r}}_0 \mathbf{E}^l = [1, 0, \dots, 0] \mathbf{E}^l, \quad (4.116)$$

meaning that

$$\lim_{l \rightarrow \infty} \bar{\mathbf{r}}_l = \mathbf{0} \text{ iff } \lim_{l \rightarrow \infty} \mathbf{E}^l = \mathbf{0}. \quad (4.117)$$

The matrix \mathbf{E} satisfies this condition if $|\lambda_{\max}| < 1$. On the other hand, if $|\lambda_{\max}| > 1$, the limit diverges. ■

Note that the theorem does not say anything about the case when $|\lambda_{\max}| = 1$.

We conclude the subsection by noting that, interestingly, weighted context-free grammars *trained* on data with maximum likelihood are *always* tight (Chi and Geman, 1998; Chi, 1999). This is not the case for some models we consider later, e.g., recurrent neural networks (cf. §5.1.2).

4.2.6 Normalizing Weighted Context-free Grammars

Having investigated probabilistic context-free grammars in terms of their tightness, we now turn out attention to general weighted context-free grammars, which define string probabilities using global normalization (cf. Eq. (4.99)). To be able to compute these probabilities, require a way to compute the normalizing constant $Z(\mathcal{G})$ and the stringsum $\mathcal{G}(\mathbf{y})$. In the section on finite-state automata, we explicitly presented an algorithm for computing the normalizing constant $Z(\mathcal{A})$. The derivation of a general allsum algorithm for weighted context-free grammars, on the other hand, is more involved and beyond the scope of this course.²⁷ Here, we simply assert that there are ways of computing the quantities in Eq. (4.99) and only consider the following result:

²⁷The allsums of individual non-terminals can be expressed as solutions to a nonlinear set of equations. Again, the interested reader should have a look at the Advanced Formal Language Theory course.

Theorem 4.2.2 (PCFGs and WCFGs are equally expressive (Smith and Johnson, 2007)). *Normalizable weighted context-free grammars with non-negative weights and tight probabilistic context-free grammars are equally expressive.*

Proof. To prove the theorem, we have to show that any WCFG can be written as a PCFG and vice versa.²⁸

\Leftarrow Since any tight probabilistic context-free grammar is simply a WCFG with $Z(\mathcal{G}) = 1$, this holds trivially.

\Rightarrow We now show that, for any WCFG, there exists a PCFG encoding the same language model. Let $\mathcal{G}_G = (\Sigma, \mathcal{N}, S, \mathcal{P}, \mathcal{W})$ be a pruned WCFG that encodes a distribution over Σ^* using Eq. (4.99). We now construct a tight probabilistic context-free grammar $\mathcal{G}_L = (\Sigma, \mathcal{N}, S, \mathcal{P}, \mathcal{W}_L)$ whose language is identical. Notice that all components of the grammar remain identical apart from the weighting function. This means that the derivations of the strings in the grammars remain the same (i.e., $\mathcal{D}_{\mathcal{G}_G} = \mathcal{D}_{\mathcal{G}_L}$)—only the weights of the derivations change, as we detail next. We define the production weights of the probabilistic CFG as follows.

$$\mathcal{W}_{\mathcal{G}_L}(X \rightarrow \alpha) \stackrel{\text{def}}{=} \frac{\mathcal{W}(X \rightarrow \alpha) \prod_{Y \in \alpha} Z(\mathcal{G}, Y)}{Z(\mathcal{G}, X)} \quad (4.118)$$

Remember that $Z(a) = 1$ for $a \in \Sigma$. Note that the assumption that \mathcal{G} is pruned means that all the quantities in the denominators are non-zero.

It is easy to see that the weight defined this way are non-negative due to the non-negativity of \mathcal{G} 's weights. Furthermore, the weights of all production rules for any non-terminal $X \in \mathcal{N}$ sum to 1, as by the definitions of $\mathcal{W}_{\mathcal{G}_L}$ and $Z(\mathcal{G}, X)$ we have

$$\sum_{X \rightarrow \alpha} \mathcal{W}_{\mathcal{G}_L}(X \rightarrow \alpha) = \sum_{X \rightarrow \alpha} \frac{\mathcal{W}(X \rightarrow \alpha) \prod_{Y \in \alpha} Z(\mathcal{G}, Y)}{Z(\mathcal{G}, X)} \quad (4.119)$$

$$= \frac{1}{Z(\mathcal{G}, X)} \sum_{X \rightarrow \alpha} \mathcal{W}(X \rightarrow \alpha) \prod_{Y \in \alpha} Z(\mathcal{G}, Y) \quad (4.120)$$

$$= \frac{1}{Z(\mathcal{G}, X)} Z(\mathcal{G}, X) \quad (4.121)$$

$$= 1 \quad (4.122)$$

We now have to show that the probabilities assigned by these two grammars match. We will do that by showing that the probabilities assigned to individual *derivations* match, implying that stringsums match as well. The probability of a derivation is defined analogously to a probability of a string, i.e., $p_{\mathcal{G}}(\mathbf{d}) = \frac{w(\mathbf{d})}{Z(\mathcal{G})}$ (where $Z(\mathcal{G}) = 1$ for tight probabilistic grammars). Let then $\mathbf{d} \in \mathcal{D}_{\mathcal{G}} = \mathcal{D}_{\mathcal{G}_L}$. Then

$$p_{\mathcal{G}_L}(\mathbf{d}) = \prod_{X \rightarrow \alpha \in \mathbf{d}} \mathcal{W}_{\mathcal{G}_L}(X \rightarrow \alpha) \quad (4.123)$$

$$= \prod_{X \rightarrow \alpha \in \mathbf{d}} \frac{\mathcal{W}(X \rightarrow \alpha) \prod_{Y \in \alpha} Z(\mathcal{G}, Y)}{Z(\mathcal{G}, X)} \quad (\text{definition of } \mathcal{W}_{\mathcal{G}_L}). \quad (4.124)$$

Notice that by multiplying over the internal nodes of the derivation tree, Eq. (4.123) includes the non-terminal allsum of each *internal* (non-root and non-leaf) non-terminal in the derivation *twice*: once as a parent of a production in the denominator, and once as a child in the numerator. These terms, therefore, all cancel out in the product. The only terms which are left are the allsums of the leaf nodes—the terminals—which are 1, and the allsum of the root node— S —which equals $Z(\mathcal{G}_G)$ and the weights of the individual productions, which multiply into the weight assigned to \mathbf{d} by the original grammar \mathcal{G}_G . This means that

$$p_{\mathcal{G}_L}(\mathbf{d}) = \frac{1}{Z(\mathcal{G}, X)} \prod_{X \rightarrow \alpha \in \mathbf{d}} \mathcal{W}(X \rightarrow \alpha) = \frac{1}{Z(\mathcal{G}, X)} w(\mathbf{d}) = p_{\mathcal{G}_G}(\mathbf{d}), \quad (4.125)$$

finishing the proof. ■

²⁸Again, by “written as”, we mean that the weighted language is the same.

This means that the classes of probabilistic and weighted context-free grammars are in fact *equally expressive*. In other words, this result is analogous to Theorem 4.1.1 in WFSAs: it shows that in the context of context-free language models, the locally normalized version of a globally-normalized model is *also* context-free.

4.2.7 Pushdown Automata

We presented context-free grammars as a formalism for specifying and representing context-free languages. Many algorithms for processing context-free languages, for example, the allsum algorithms and their generalizations, can also be directly applied to context-free grammars. However, it is also convenient to talk about processing context-free languages in terms of computational models in the form of automata, i.e., the *recognizer* of the language.²⁹ As we mentioned, the types of automata we considered so far, (weighted) finite-state automata, can only recognize regular languages. To recognize context-free languages, we must therefore extend finite-state automata.³⁰ We do that by introducing **pushdown automata** (PDA), a more general and more expressive type of automata.

Single-stack Pushdown Automata

Pushdown automata augment finite-state automata by implementing an additional *stack* memory structure for storing *arbitrarily long* strings from a designated alphabet, which allows them to work with unbounded memory effectively. Abstractly, this unbounded memory is the only difference to finite-state automata. However, the definition looks slightly different:

Definition 4.2.31. A **pushdown automaton** (PDA) is a tuple $\mathcal{P} \stackrel{\text{def}}{=} (\Sigma, Q, \Gamma, \delta, (q_i, \gamma_i), (q_f, \gamma_f))$, where:

- Q is a finite set of states;
- Σ is a finite set of input symbols called the input alphabet;
- Γ is a finite set of stack symbols called the stack alphabet;
- $\delta \subseteq Q \times \Gamma^* \times (\Sigma \cup \{\varepsilon\}) \times Q \times \Gamma^*$ is a multiset representing the transition function;
- (q_i, γ_i) is called the initial configuration and (q_f, γ_f) is called the final configuration, where $q_i, q_f \in Q$ and $\gamma_i, \gamma_f \in \Gamma^*$.

The initial and final configurations in pushdown play analogous roles to the sets of initial and final sets of finite-state automata. Compared to the latter, they also allow for different starting configurations of the stack coupled with each possible initial or final state.

Stacks are represented as strings over Γ , from bottom to top. Thus, in the stack $\gamma = X_1 X_2 \cdots X_n$, the symbol X_1 is at the bottom of the stack, while X_n is at the top. $\gamma = \emptyset$ denotes the empty stack.

Definition 4.2.32. A **configuration** of a PDA is a pair (q, γ) , where $q \in Q$ is the current state and $\gamma \in \Gamma^*$ is the current contents of the stack.

The initial and final configurations of a PDA are examples of configurations; it is possible to generalize the initial and final stacks to (say) regular expressions over Γ , but the above definition suffices for our purposes.

A PDA moves from configuration to configuration by following transitions of the form $q \xrightarrow{a, \gamma_1 \rightarrow \gamma_2} r$, which represents a move from the state q to state r , while popping the sequence of symbols $\gamma_1 \in \Gamma^*$ from the top of the stack and pushing the sequence $\gamma_2 \in \Gamma^*$. The PDA transition function therefore not only depends on the current state q and input symbol a , but also on some *finite* sequence of symbols on the *top* of the stack. The stack hence determines the behavior of the automaton, and since the set of possible configurations of the stack is infinite, the set of configurations of the automaton is infinite, in contrast to finite-state automata.

To describe how pushdown automata process strings, we introduce the concepts of scanning and runs.

²⁹This relationship between a formalism specifying how to *generate* (i.e., a grammar) and a model of *recognizing* a language can be seen in multiple levels of the hierarchy of formal languages. In the case of context-free languages, the former are context-free grammars, while the latter are pushdown automata discussed in this subsection. Regular languages as introduced in the previous section, however, are simply defined in terms of their recognizers—finite-state automata.

³⁰Formally, we would of course have to prove that finite-state automata cannot model context-free languages. This can be done with the so-called pumping lemma, which are outside the scope of this class.

Definition 4.2.33. We say that $\tau = (p, \gamma_1, a, q, \gamma_2) \in \delta$ **scans** a , and if $a \neq \varepsilon$, we call τ **scanning**; otherwise, we call it **non-scanning**.

Definition 4.2.34. If $(q_1, \gamma\gamma_1)$ and $(q_2, \gamma\gamma_2)$ are configurations, and τ is a transition $q_1 \xrightarrow{a, \gamma_1 \rightarrow \gamma_2} q_2$, we write $(q_1, \gamma\gamma_1) \Rightarrow_\tau (q_2, \gamma\gamma_2)$.

Since the behavior of a pushdown automaton does not only depend on the states encountered by it but also on the content of the stack, we generalize the notion of a path to include the *configuration* of the automaton. This is called a run.

Definition 4.2.35. A **run** of a PDA \mathcal{P} is a sequence of configurations and transitions

$$\pi = (q_0, \gamma_0), \tau_1, (q_1, \gamma_1), \dots, \tau_n, (q_N, \gamma_N)$$

where, for $n = 1, \dots, N$, we have $(q_{n-1}, \gamma_{n-1}) \Rightarrow_{\tau_n} (q_n, \gamma_n)$.³¹ A run is called **accepting** if (q_0, γ_0) is the initial configuration and (q_N, γ_N) is the final configuration. If, for $n = 1, \dots, N$, τ_n scans a_n , then we say that π scans the string $a_1 \cdots a_N$. We write $\Pi(\mathcal{P}, \mathbf{y})$ for the set of runs that scan \mathbf{y} and $\Pi(\mathcal{P})$ for the set of all accepting runs of \mathcal{P} .

Definition 4.2.36. We say that the PDA \mathcal{P} **recognizes** the string \mathbf{y} if $\Pi(\mathcal{P}, \mathbf{y}) \neq \emptyset$, i.e., if there exists an accepting run with the yield \mathbf{y} . The set of all strings recognized by \mathcal{P} is the **language recognized by \mathcal{P}** , which we denote by $L(\mathcal{P})$, i.e.,

$$L(\mathcal{P}) \stackrel{\text{def}}{=} \{\mathbf{y} \mid \Pi(\mathcal{P}, \mathbf{y}) \neq \emptyset\}. \quad (4.126)$$

Example 4.2.12. Fig. 4.8 shows an example of a pushdown automaton \mathcal{P} accepting the language $L(\mathcal{P}) = \{a^n b^n \mid n \in \mathbb{N}\}$. $(1 \xrightarrow{a, \varepsilon \rightarrow X} 1, 1 \xrightarrow{\varepsilon, \varepsilon \rightarrow \varepsilon} 2)$ is a run of \mathcal{P} ; $(1 \xrightarrow{a, \varepsilon \rightarrow X} 1, 1 \xrightarrow{\varepsilon, \varepsilon \rightarrow \varepsilon} 2, 2 \xrightarrow{b, X \rightarrow \varepsilon} 2)$ is an accepting run of \mathcal{P} .

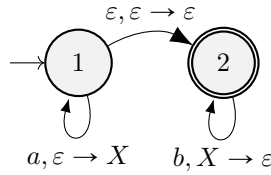


Figure 4.8: The PDA that accepts the language $\{a^n b^n \mid n \in \mathbb{N}\}$.

Lastly, we define deterministic pushdown automata, analogously to their finite-state version (Definition 4.1.3). Recall that in the case of finite-state automata, a deterministic machine has at most one possible next move for each state. Similarly, a deterministic pushdown automaton has at most one possible next move for each *configuration*.

Definition 4.2.37. A PDA $\mathcal{P} = (\Sigma, Q, \Gamma, \delta, (q_i, \gamma_i), (q_f, \gamma_f))$ is **deterministic** if

- there are no transitions of the type $(q, \varepsilon, \gamma, p, \gamma)$;
- for every $(q, a, \gamma) \in Q \times \Sigma \cup \{\varepsilon\} \times \Gamma^*$, there is at most one transition $(q, a, \gamma, p, \gamma') \in \delta$;
- if there is a transition $(q, a, \gamma, p, \gamma') \in \delta$ for some $a \in \Sigma$, then there is no transition $(q, \varepsilon, \gamma, p, \gamma'') \in \delta$.

Otherwise, \mathcal{P} is **non-deterministic**.

Importantly, *not all* context-free languages can be recognized by deterministic pushdown automata. That is, in contrast to finite-state automata, where deterministic machines are just as powerful as non-deterministic ones (at least in the unweighted case—interestingly, some weighted non-deterministic FSAs cannot be determinized), non-deterministic pushdown automata are more expressive than deterministic ones. Specifically, as stated in Theorem 4.2.3, non-deterministic pushdown automata recognize exactly context-free languages, while deterministic pushdown automata only recognize a subset of them (Sipser, 2013).

³¹Sometimes it will be convenient to treat π as a sequence of only configurations or only transitions.

Weighted Pushdown Automata

Analogously to the finite-state case, and the case of context-free grammars, we now also extend the definition of a pushdown automaton to the weighted case. The formal definition is:

Definition 4.2.38. A *real-weighted pushdown automaton* (WPDA) is a tuple $\mathcal{P} = (Q, \Sigma, \Gamma, \delta, (q_i, \gamma_i), (q_f, \gamma_f))$, where:

- Q is a finite set of states;
- Σ is a finite set of input symbols called the input alphabet;
- Γ is a finite set of stack symbols called the stack alphabet;
- $\delta \subseteq Q \times \Gamma^* \times (\Sigma \cup \{\varepsilon\}) \times Q \times \Gamma^* \times \mathbb{R}$ is a multi-set representing the transition weighting function;
- (q_i, γ_i) is called the initial configuration and (q_f, γ_f) is called the final configuration, where $q_i, q_f \in Q$ and $\gamma_i, \gamma_f \in \Gamma^*$.

As you can see, the only difference between the weighted and the unweighted case is the transition function, which in the weighted case weights the individual transitions instead of specifying the set of possible target configurations.

As with WFSAs (Definition 4.1.17) and WCFGs (Definition 4.2.16), we now define probabilistic WPDA. This definition, however, is a bit more subtle. Notice that the transition weighting “function” δ in a WPDA is crucially still a *finite*—there is only a finite number of actions we can ever do. Similarly, when defining a *probabilistic* PDA, we have to limit ourselves to a finite number of configurations over which we define probability distributions over the next possible actions.

We define a probabilistic pushdown automaton given an equivalence relation as follows.

Definition 4.2.39. A WPDA $\mathcal{P} = (Q, \Sigma, \Gamma, \delta, (q_i, \gamma_i), (q_f, \gamma_f))$ is **probabilistic** if it holds that

$$\forall q \frac{a, \gamma_1 \rightarrow \gamma_2 / w}{r \in \delta} : w \geq 0 \quad (4.127)$$

and for any $q \in Q$ and $\gamma \in \Gamma^*$

$$\sum_{\substack{q \xrightarrow{a, \gamma_1 \rightarrow \gamma_2 / w} r \\ \text{s.t. } \gamma_1 \triangleleft \gamma}} w = 1. \quad (4.128)$$

Definition 4.2.40. If (q_1, γ_1) and (q_2, γ_2) are configurations, and τ is a transition $q_1 \xrightarrow{a, \gamma_1 \rightarrow \gamma_2 / w} q_2$ with $w \neq 0$, we write $(q_1, \gamma_1) \Rightarrow_\tau (q_2, \gamma_2)$.

Definition 4.2.41. If $\delta(p, \gamma_1, a, q, \gamma_2) = w$, then we usually write

$$\tau(p, \gamma_1 \xrightarrow{a} q, \gamma_2) = w \quad (4.129)$$

or that δ has transition $(q \xrightarrow{a, \gamma_1 \rightarrow \gamma_2 / w} p)$. We sometimes let τ stand for a transition, and we define $\delta(\tau) = w$.

And again, just like we combined the weights of individual transitions into the weights of paths in WFSAs, and we combined the weights of production rules into the weights of the trees in WCFGs, we now multiplicatively combine the weights of individual transitions in a run to define the weight of a run in a WPDA:

Definition 4.2.42. The **weight** $w(\pi)$ of a run

$$\pi = (q_0, \gamma_0), \tau_1, (q_1, \gamma_1), \dots, \tau_N, (q_N, \gamma_N)$$

is the multiplication of the transition weights, i.e.,

$$w(\pi) \stackrel{\text{def}}{=} \prod_{n=1}^N \delta(\tau_n) \quad (4.130)$$

Analogously to a stringsum in WFSAs, we define the stringsum for a string \mathbf{y} in a WPDA \mathcal{P} as the sum over the weights of all runs scanning \mathbf{y} .

Definition 4.2.43. Let \mathcal{P} be a WPDA and $\mathbf{y} \in \Sigma^*$ a string. The **stringsum** for \mathbf{y} in \mathcal{P} is defined as

$$\mathcal{P}(\mathbf{y}) \stackrel{\text{def}}{=} \sum_{\pi \in \Pi(\mathcal{P}, \mathbf{y})} w(\pi) \quad (4.131)$$

Definition 4.2.44. We say that the PDA \mathcal{P} **recognizes** the string \mathbf{y} with the weight $\mathcal{P}(\mathbf{y})$.

With this, we can define the weighted language defined by a WPDA.

Definition 4.2.45. Let \mathcal{P} be a WPDA. The **(weighted) language** $L(\mathcal{P})$ of \mathcal{P} is defined as

$$L(\mathcal{P}) \stackrel{\text{def}}{=} \{(\mathbf{y}, \mathcal{P}(\mathbf{y})) \mid \mathbf{y} \in \Sigma^*\} \quad (4.132)$$

Finally, we also define the WPDA allsum and normalizable WPDAs.

Definition 4.2.46 (WPDA Allsum). The **allsum** of a WPDA \mathcal{P} is defined as

$$Z(\mathcal{P}) \stackrel{\text{def}}{=} \sum_{\pi \in \Pi(\mathcal{P})} w(\pi) \quad (4.133)$$

Definition 4.2.47. A WPDA \mathcal{P} is **normalizable** if $Z(\mathcal{P})$ is finite, i.e., if $Z(\mathcal{P}) < \infty$.

Relationship to Context-free Grammars

We motivated the introduction of pushdown automata as a means of recognizing context-free languages. However, this correspondence is not obvious from the definition! Indeed, the equivalence of the expressive power of context-free grammars and pushdown automata is a classic result in formal language theory, and it is summarised by the theorem below:

Theorem 4.2.3. A language is context-free if and only if some pushdown automaton recognizes it.

Proof. See Theorem 2.20 in Sipser (2013). ■

This result extends to the probabilistic case.

Theorem 4.2.4. A language is generated by a probabilistic context-free grammar if and only if some probabilistic pushdown automaton recognizes it.

Proof. See Theorems 3 and 7 in Abney et al. (1999). ■

Lastly, analogously to how Theorem 4.2.2 showed that weighted context-free grammars are equally expressive as probabilistic context-free grammars, the following theorem asserts the same about pushdown automata:

Theorem 4.2.5. Any globally normalized weighted pushdown automaton can be locally normalized. More precisely, this means the following. Let \mathcal{P} be a weighted pushdown automaton. Then, there exists a probabilistic pushdown automaton \mathcal{P}_p such that

$$\mathcal{P}_p(\mathbf{y}) = \frac{\mathcal{P}(\mathbf{y})}{Z(\mathcal{P})} \quad (4.134)$$

for all $\mathbf{y} \in \Sigma^*$.

Proof. The proof is not straightforward: it can be shown that one cannot simply convert an arbitrary weighted pushdown automaton into a locally-normalized one directly. Rather, the construction of the latter goes through their *context-free grammars*: given a WPDA \mathcal{P} , one first constructs the WCFG equivalent to \mathcal{P} , and then converts that to a locally normalized one (cf. Theorem 4.2.2), i.e., a PCFG. Then, this PCFG can be converted to a structurally quite different probabilistic pushdown automaton \mathcal{P}_p , which nevertheless results in the language we require (Eq. (4.134)).

This construction is described in more detail in Abney et al. (1999) and Butoi et al. (2022). ■

4.2.8 Pushdown Language Models

We can now define pushdown language models, the title of this section.

Definition 4.2.48. A *pushdown language model* is a language model whose weighted language equals the language of some weighted pushdown automaton, i.e., if there exists a weighted pushdown automaton \mathcal{P} such that $L(\mathcal{P}) = L(p_{LM})$.

Similarly, pushdown automata also induce language models.

Definition 4.2.49. Let \mathcal{P} be a weighted pushdown automaton. We define the *language model induced by \mathcal{P}* as the probability distribution induced by the probability mass function

$$p_{LM\mathcal{P}}(\mathbf{y}) \stackrel{\text{def}}{=} \frac{\mathcal{P}(\mathbf{y})}{Z(\mathcal{P})}, \quad (4.135)$$

for any $\mathbf{y} \in \Sigma^*$.

You might wonder why we specifically define pushdown language models and models induced by them if WPDAs are equivalent to WCFGs (cf. §4.2.7). In that sense, a language model is context-free if and only if it is a pushdown language model. However, this holds only for single-stack pushdown automata which we have discussed so far. We make this explicit distinction of pushdown language models with an eye to the next section, in which we introduce *multi-stack* WPDAs. Those are, as it turns out, much more powerful (expressive) than context-free grammars. We will, however, reuse this definition of a pushdown language model for those more powerful machines.

4.2.9 Multi-stack Pushdown Automata

We now consider an extension of (weighted) pushdown automata, namely, machines that employ *multiple* stacks. While this might not seem like an important distinction, we will see shortly that this augmentation results in a big difference in the expressiveness of the framework!

Definition 4.2.50. A *two-stack pushdown automaton* (2-PDA) is a tuple $\mathcal{P} = (\Sigma, Q, \Gamma_1, \Gamma_2, \delta, (q_\iota, \gamma_{\iota 1}, \gamma_{\iota 2}), (q_\varphi, \gamma_{\varphi 1}, \gamma_{\varphi 2}))$ where:

- Σ is a finite set of input symbols called the input alphabet;
- Q is a finite set of states;
- Γ_1 and Γ_2 are finite sets of stack symbols called the stack alphabets;
- $\delta \subseteq Q \times \Gamma_1^* \times \Gamma_2^* \times (\Sigma \cup \{\varepsilon\}) \times Q \times \Gamma_1^* \times \Gamma_2^*$ is a multiset representing the transition function;
- $(q_\iota, \gamma_{\iota 1}, \gamma_{\iota 2})$ is called the initial configuration and $(q_\varphi, \gamma_{\varphi 1}, \gamma_{\varphi 2})$ is called the final configuration, where $q_\iota, q_\varphi \in Q$, $\gamma_{\iota 1}, \gamma_{\varphi 1} \in \Gamma_1^*$, and $\gamma_{\iota 2}, \gamma_{\varphi 2} \in \Gamma_2^*$.

Note that we could more generally define a k -stack PDA by including k stacks in the definition, but the restriction to two stacks will be sufficient for our needs, as we will see in the next subsection. The transition function now depends on the values stored in *both* of the stacks. The definitions of the configuration and run of a two-stack PDA are analogous to the single-stack variant, with the addition of the two stacks. We again extend this definition to the weighted and the probabilistic case.

Definition 4.2.51. A *two-stack real-weighted pushdown automaton* (2-WPDA) is a tuple $\mathcal{P} = (\Sigma, Q, \Gamma_1, \Gamma_2, \delta, (q_\iota, \gamma_{\iota 1}, \gamma_{\iota 2}), (q_\varphi, \gamma_{\varphi 1}, \gamma_{\varphi 2}))$, where:

- Σ is a finite set of input symbols called the input alphabet;
- Q is a finite set of states;
- Γ_1 and Γ_2 are finite sets of stack symbols called the stack alphabets;

- $\delta \subseteq Q \times \Gamma_1^* \times \Gamma_2^* \times (\Sigma \cup \{\varepsilon\}) \times Q \times \Gamma_1^* \times \Gamma_2^* \times \mathbb{R}$ is a multiset representing the transition weighting function;
- $(q_i, \gamma_{i1}, \gamma_{i2})$ is called the initial configuration and $(q_f, \gamma_{f1}, \gamma_{f2})$ is called the final configuration, where $q_i, q_f \in Q$, $\gamma_{i1}, \gamma_{f1} \in \Gamma_1^*$, and $\gamma_{i2}, \gamma_{f2} \in \Gamma_2^*$.

And lastly, we define probabilistic two-stack PDAs:

Definition 4.2.52. A 2-WPDA $\mathcal{P} = (Q, \Sigma, \Gamma_1, \Gamma_2, \delta, (q_i, \gamma_{i1}, \gamma_{i2}), (q_f, \gamma_{f1}, \gamma_{f2}))$ is **probabilistic** if for any configuration (q, γ_1, γ_2) it holds that

$$\forall q \xrightarrow{a, \gamma_1 \rightarrow \gamma'_1, \gamma_2 \rightarrow \gamma'_2 / w} r \in \delta : w \geq 0 \quad (4.136)$$

and for any $q \in Q$ and $\gamma \in \Gamma_1^*$, $\gamma' \in \Gamma_2^*$

$$\sum_{\substack{q \xrightarrow{a, \gamma_1 \rightarrow \gamma'_1, \gamma_2 \rightarrow \gamma'_2 / w} r \\ \text{s.t. } \gamma_1 \triangleleft \gamma \text{ and } \gamma_2 \triangleleft \gamma'}} w = 1. \quad (4.137)$$

Turing Completeness of Multi-stack Pushdown Automata

Besides modeling more complex languages than finite-state language models from §4.1, (multi-stack) pushdown automata will also serve an important part in analyzing some modern language models that we introduce later. Namely, we will show that Recurrent neural networks (cf. §5.1.2) can *simulate* any two-stack PDA. This will be useful when reasoning about the computational expressiveness of recurrent neural networks because of a fundamental result in the theory of computation, namely, that two-stack PDAs are *Turing complete*:

Theorem 4.2.6. Any 2-stack pushdown automaton is Turing complete.

Proof. The equivalence is quite intuitive: the two stacks (which are infinite in one direction) of the 2-PDA can simulate the tape of a Turing machine by popping symbols from one stack and pushing symbols onto the other one simultaneously. The head of the Turing machine then effectively reads the entries at the top of one of the two stacks. For a formal proof, see Theorem 8.13 in Hopcroft et al. (2006). ■

This is also the reason why we only have to consider two-stack PDAs—they can compute everything that can be computed, meaning that additional stacks do not increase their expressiveness!

Since unweighted pushdown automata are simply special cases of weighted PDAs, which are equivalent to probabilistic PDAs, we can therefore also conclude:

Corollary 4.2.1. Any weighted 2-stack pushdown automaton is Turing complete.

Corollary 4.2.2. Any probabilistic 2-stack pushdown automaton is Turing complete.

A straightforward consequence of the Turing completeness of two-stack PPDA is that their *tightness* is undecidable.

Theorem 4.2.7 (Tightness of 2-PPDA is undecidable). *The tightness of a probabilistic two-stack pushdown automaton is undecidable.*

Proof. We start with a simple observation: a pushdown automaton \mathcal{P} is tight if and only if it *halts* on inputs with measure 1 (given the probability measure on $\Sigma^* \cup \Sigma^\infty$ defined in §2.5.4), as this, by the definition of the language accepted by the WPDA (cf. Definition 4.2.44), corresponds to its language only containing *finite* strings with probability 1.

Let \mathcal{M} then be a Turing machine and \mathcal{P} be a 2-PPDA which simulates it. Then, \mathcal{P} is tight if and only if it halts with probability 1 (again, based on the probability measure from above). This is equivalent to the problem of \mathcal{M} halting with probability 1—this, however, is a variant of the **halting problem**, which is one of the fundamental undecidable problems. We have therefore reduced the problem of determining the tightness of 2-PPDAs to the halting problem, implying that the former is undecidable. ■

You might wonder what this means for the (weighted) languages recognized by multiple-stack (weighted) automata. Turing machines can recognize *recursively enumerable languages*. This means that weighted multi-stack pushdown automata model *distributions* over recursively enumerable languages. To see why this might be useful, let us finish the discussion of context-free languages with an example of a language model that is *not* context-free:

Example 4.2.13. Let $\Sigma = \{a\}$ and $p_{LM}(a^n) = e^{-\lambda} \frac{\lambda^n}{n!}$ for $n \in \mathbb{N}_{\geq 0}$, i.e., $L(p_{LM}) \ni \mathbf{y} = a^n \sim \text{Poisson}(\lambda)$ for some $\lambda > 0$. This language is not context-free: the proof, however, is not trivial. We direct the reader to [Icard \(2020b\)](#) for one.

4.3 Exercises

Exercise 4.1

Prove the following lemma.

Lemma 4.3.1. *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ and $q \in Q$. Then*

$$Z(\mathcal{A}, q) = \sum_{q \xrightarrow{a/w} q' \in \delta_{\mathcal{A}_L}} \omega\left(q \xrightarrow{a/\cdot} q'\right) Z(\mathcal{A}, q') + \rho(q) \quad (4.138)$$

Exercise 4.2

Show that the expression for the log-likelihood of the n -gram model can be rewritten as

$$\ell\ell(\mathcal{D}) = \sum_{m=1}^M \sum_{t=1}^{|\mathbf{y}^{(m)}|} \log \theta_{y_n | \mathbf{y} < n} = \sum_{\substack{\mathbf{y} \\ |\mathbf{y}|=n}} C(\mathbf{y}) \theta_{y_n | \mathbf{y} < n} \quad (4.139)$$

with the quantities as defined in Proposition 4.1.3. This is a common trick. It is also known as the **token to type switch** because we switch from counting over the individual tokens to counting over their identities (types)

Exercise 4.3

Let $C(\mathbf{y})$ be the string occurrence count for $\mathbf{y} \in \Sigma^*$ occurrence count as defined in Proposition 4.1.3. Show (or simply convince yourself) that, in a given training corpus \mathcal{D}

$$\sum_{y' \in \Sigma} C(y_1 \dots y_{n-1} y') = C(y_1 \dots y_{n-1}) \quad (4.140)$$

Chapter 5

Neural Network Language Models

Chapter 4 introduced two classical language modeling frameworks: finite-state language models and context-free language models. While those served as a useful introduction to the world of language modeling, most of today’s state-of-the-art language models go beyond the modeling assumptions of these two frameworks. This chapter dives into the diverse world of modern language modeling architectures, which are based on neural networks. We define two of the most common architectures—recurrent neural networks and transformers—and some of their variants. The focus is again on rigorous formalization and theoretical understanding—we analyze the introduced models in terms of the theoretical foundations so far (e.g., expressiveness and tightness)—but, due to their practical applicability, we also study some practical aspects of the models.

We begin with recurrent neural networks.

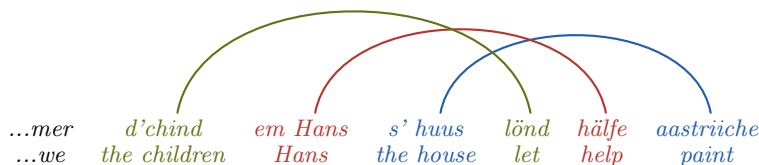
5.1 Recurrent Neural Language Models

The first neural language modeling architecture we consider is one based on recurrent neural networks. Recurrent neural networks capture the idea of the *sequential* processing of strings relatively naturally while also making decisions based on an *infinite* context. Before delving into the technical details of recurrent neural networks (RNNs), however, we first motivate the introduction of modeling contexts of unbounded length. Then, we formally define recurrent neural networks and devote a large portion of the section to their theoretical properties. The most important of those will be the Turing completeness of this architecture, as it has numerous consequences on the solvability of many of the tasks we might be interested in, such as finding the most probable string in the language model represented by a recurrent neural network and determining whether an RNN is tight.

5.1.1 Human Language is Not Context-free

Recall that we motivated the introduction of context-free languages by observing that finite memory is insufficient to model all formal phenomena of human language, e.g., infinite recursion (cf. Example 4.2.1). Context-free languages described by context-free grammars and pushdown automata were able to capture those. However, human language is more expressive than that—it includes linguistic phenomena that cannot be described by context-free grammars. A typical example is called **cross-serial dependencies**, which are common in *Swiss German*.

Example 5.1.1 (Cross-Serial Dependencies in Swiss German, [Shieber, 1985](#)). *Swiss German is a textbook example of a language with grammatical cross-serial dependencies, i.e., dependencies in which the arcs representing them, cross. In the example sentence below, the words connected with arcs are objects and verbs belonging to the same predicates (verb phrases). Because of that, they have to agree on the form—they depend on one another. As we show next, context-free languages cannot capture such dependencies.*



Why are cross-serial dependencies non-context-free? Before reasoning about the phenomenon of cross-serial dependencies, we revisit Example 4.2.1 with a somewhat more formal approach. The arbitrarily deep nesting can, for example, be abstractly represented with the expression

$$xA^nB^ny \quad (5.1)$$

with¹

$$\begin{aligned} x &= \text{"The cat"} \\ A &= \text{"the dog"} \\ B &= \text{"barked at"} \\ y &= \text{"likes to cuddle"}. \end{aligned}$$

From this abstract perspective, center embeddings are very similar to the D(1) language (Example 4.2.5), in that every noun phrase “the dog” has to be paired with a verb phrase “barked at”, which cannot be represented by any regular language.

In a similar fashion, Example 5.1.1 can abstractly be represented with the expression

$$xA^mB^nC^myD^nz \quad (5.2)$$

¹In this case, we of course only consider an arbitrarily long sequence of barking dogs.

with

$$\begin{aligned} x &= \text{"...mer"} \\ A &= \text{"d'chind"} \\ B &= \text{"em Hans"} \\ y &= \text{"s' huus"} \\ C &= \text{"lönd"} \\ D &= \text{"hälfe"} \\ z &= \text{"aastriiche"}. \end{aligned}$$

Admittedly, this is a relatively uncommon formulation even with $n = m = 1$. It should be taken with a grain of salt, as the title of the original publication discussing this phenomenon, *Evidence against context-freeness* (Shieber, 1985), also suggests. However, theoretically, the number of repetitions of “d’chind” and “lönd”, as well as “em Hans” and “hälfe”, can be increased arbitrarily. Repeating the former would correspond to having many groups of children. The last of the groups would let Hans help paint the house, whereas each of the previous groups would let the group after them either let Hans paint the house or recurse onto another group of children. Similarly, repeating “em Hans” and “hälfe” would correspond to a number of Hanses, each either helping another Hans or helping paint the house. Then, using the pumping lemma for *context-free* languages, it can be shown that the expressions of the form in Eq. (5.2) cannot be recognized by any context-free grammar. We refer the readers to Hopcroft et al. (2006, Example 7.20) for detailed proof.

Example 5.1.1 means that to model a human language formally, we need more expressive formalisms than context-free grammars or pushdown automata as described in the previous sections.² However, instead of defining a more expressive formalism motivated by formal language theory (like we did with context-free grammars and center embeddings), we now introduce recurrent neural networks, which, as we will see, under certain assumptions, have the capacity to model all computable languages (i.e., they are Turing complete). Moreover, they can also model *infinite* lengths of the context $\mathbf{y}_{<t}$ in a very flexible way. In the next section, we define them formally.

5.1.2 Recurrent Neural Networks

As discussed, natural languages are beyond the descriptive power of regular and context-free languages. Now, we turn to a class of models that is theoretically capable of recognizing all computable languages: **recurrent neural networks** (RNNs).³

An Informal Introduction to Recurrent Neural Networks

Human language is inherently sequential: we produce and consume both spoken as well as written language as a stream of units.⁴ This structure is reflected in some of the algorithms for processing language we have seen so far. For example, finite-state automata (cf. §4.1) process the input string one symbol at a time and build the representations of the string seen so far in the current state of the automaton. Pushdown automata function similarly, but additionally keep the stack as part of the configuration.

Recurrent neural networks are neural networks that capture the same idea of iterative processing of the input but do so in a more flexible way than the finite-memory finite-state automata and the stack-based pushdown automata. Very abstractly, a recurrent neural network sequentially processes a sequence of inputs and, while doing so, produces a sequence of **hidden states**, which we will denote as \mathbf{h} , based on a transition function in form of a **recurrent dynamics map**, which acts similarly to a (deterministic) transition function in a finite-state machine: given the current hidden state and an input symbol, it (deterministically) determines

²On the other hand, note that we would ideally also like to upper-bound the expressive power of the formal models, as this introduces useful inductive biases for learning and sparks insights into how humans process language. This means that we would not simply like to jump to Turing-complete models in such an exploration of language models.

³In this subsection, we focus on the applications of recurrent neural networks to language modeling. However, recurrent neural networks have been widely used to process sequential data and time series, thanks to their power of taking in arbitrary-length inputs.

⁴So far, we have simply referred to those units as *symbols*.

the next hidden state. The hidden states play, as we will see, an analogous role to the states of a finite-state automaton or the configuration of a pushdown automaton: The current hidden state of a recurrent neural network at time t determines, together with the input at time t , through the dynamics map, the hidden state at time $t + 1$ —indeed, very similar to how finite-state automata process strings and transition between their states. Again, the hidden state can be thought of as a compact (constant-size) summary of the input $\mathbf{y}_{\leq t}$ seen so far and should ideally characterize $\mathbf{y}_{\leq t}$ as well as possible (in the sense of retaining all information required for continuing the string). Remember from §4.2.1 that the finite number of states of a finite-state automaton presented a serious limitation to its ability to model human language. As we will see, the main difference between (weighted) finite-state automata and RNNs is that the latter can work with *infinite* state spaces, for example, \mathbb{R}^D in the abstract formulation, or \mathbb{Q}^D in a digital computing system, such as a computer. This, together with the flexibility of the transition function between hidden states, will allow RNNs to represent more complex languages than those recognized by finite-state automata or context-free grammars. In fact, the large state space and the flexible transition functions endow RNNs, under some assumptions, with the possibility to model infinitely-long-term dependencies on the input string, distinguishing them from the Markovian n -gram models.

You might wonder why we refer to the current state of an RNN as *hidden* states instead of only as *states*, as with finite-state automata. Indeed, when analyzing recurrent neural networks in terms of their expressivity and connections to classical models of computation, we will regard the hidden states as completely analogous to states in a finite-state or pushdown automaton. The hidden part comes from the fact that the hidden states \mathbf{h} are usually not what we are interested in when modeling language with an RNN. Rather, \mathbf{h} is simply seen as a component in a system that produces individual *conditional probabilities* over the next symbol, as in sequence models (cf. Definition 2.5.2)—these conditional probabilities are the actual “visible” parts, while the “internal” states are, therefore, referred to as hidden.

RNNs are abstractly illustrated in different ways in the literature. Often, they are represented as a sequence of hidden states and the input symbols consumed to arrive at those states—this is shown in Fig. 5.1a. They can also be presented more similarly to automata, with (a possibly infinite) labeled graph, where the transition labels again correspond to the symbols used to enter the individual states. This is presented in Fig. 5.1b. Lastly, due to the infinite state space, one can also think of an RNN as a system that keeps the most current hidden state in memory and updates it as new symbols are consumed—this is shown in Fig. 5.1c.⁵

A Formal Definition of Recurrent Neural Networks

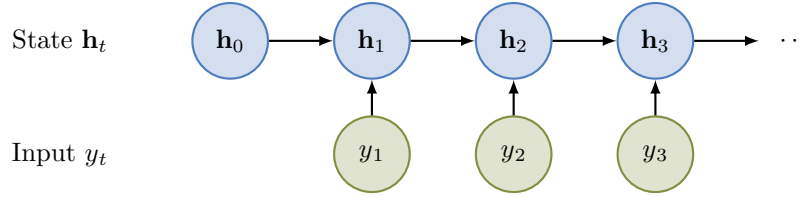
Having introduced RNNs and their motivations informally, we now move to their formal definition. Our definition and treatment of recurrent neural networks might differ slightly from what you might normally encounter in the literature. Namely, we define RNNs below as abstract systems transitioning between possibly infinitely-many states. Our definition will allow for an intuitive connection to classical language models such as finite-state and pushdown language models as well as for tractable theoretical analysis in some special cases. Specifically, when analyzing RNNs theoretically, we will make use of their connections to automata we saw in Chapter 4.

In an abstract sense, recurrent neural networks can be defined as a system transitioning between possibly infinitely-many states, which we will assume to be vectors in a vector space. Specifically, we will distinguish between *real* and *rational* recurrent neural networks.

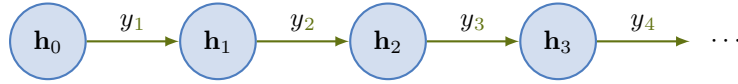
Definition 5.1.1 (Real-valued Recurrent Neural Network). *Let Σ be an alphabet. A (deterministic) **real-valued recurrent neural network** \mathcal{R} is a four-tuple $(\Sigma, D, \mathbf{f}, \mathbf{h}_0)$ where*

- Σ is the alphabet of input symbols;
- D is the dimension of \mathcal{R} ;
- $\mathbf{f}: \mathbb{R}^D \times \Sigma \rightarrow \mathbb{R}^D$ is the dynamics map, i.e., a function defining the transitions between subsequent states;

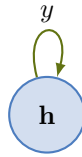
⁵More precisely, these illustrations correspond to *first-order* RNNs, which are by far the most common. Later, we will also briefly consider higher-order RNNs, whose hidden state update depends on multiple previous hidden states.



(a) An abstract depiction of how an RNN processes one symbol in a string. The hidden state \mathbf{h}_t summarizes the inputs $y_1 y_2 \dots y_t$.



(b) An abstract depiction of an RNN as an automaton. The transitions between the possibly infinitely-many hidden states are determined by the dynamics map.



(c) An abstract depiction of an RNN as a system updating the hidden state \mathbf{h} depending on the input y .

Figure 5.1: Different possible depictions of an abstract RNN model. The way that the hidden states are updated based on the input symbol y_t is abstracted away.

- $\mathbf{h}_0 \in \mathbb{R}^D$ is the *initial state*.

We analogously define **rational-valued recurrent neural networks** as recurrent neural networks with the hidden state space \mathbb{Q}^D instead of \mathbb{R}^D . You might wonder why we make the distinction. Soon, when we take on theoretical analysis of RNNs, it will become important over which state spaces the models are defined. RNNs implemented in a computer using floating-point numbers, of course, cannot have irrational-valued weights—in this sense, all implemented recurrent neural networks are rational. However, defining the models over the real numbers crucially allows us to perform operations from calculus for which some sort of continuity and smoothness is required, for example, differentiation for gradient-based learning (cf. §3.2.3).

Example 5.1.2. *An example of a rational-valued RNN is the series*

$$h_t = \frac{1}{2}h_{t-1} + \frac{1}{h_{t-1}} \quad (5.3)$$

which we considered in Example 3.1.1. In this case

- $\Sigma = \{a\}$
- $D = 1$
- $\mathbf{f}: (x, a) \mapsto \frac{1}{2}x + \frac{1}{x}$
- $\mathbf{h}_0 = 2$

Example 5.1.3. *The tuple $\mathcal{R} = (\Sigma, D, \mathbf{f}, \mathbf{h}_0)$ where*

- $\Sigma = \{a, b\}$
- $D = 2$

- $\mathbf{f}: (\mathbf{x}, y) \mapsto \begin{cases} \begin{pmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{pmatrix} \mathbf{x} & \text{if } y = a \\ \begin{pmatrix} \cos \psi & -\sin \psi \\ \sin \psi & \cos \psi \end{pmatrix} \mathbf{x} & \text{otherwise} \end{cases}$
- $\mathbf{h}_0 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

is an example of a real-valued RNN which rotates the current hidden state by the angle ϕ if the input symbol is a and rotates it by ψ if the symbol is b .

Example 5.1.4. Another example of an RNN would be the tuple

- $\Sigma = \text{GOOD} \cup \text{BAD} = \{\text{"great"}, \text{"nice"}, \text{"good"}\} \cup \{\text{"awful"}, \text{"bad"}, \text{"abysmal"}\}$
- $D = 2$
- $\mathbf{f}: (\mathbf{h}, a) \mapsto \begin{cases} \mathbf{h} + \begin{pmatrix} 1 \\ 0 \end{pmatrix} & \text{if } a \in \text{GOOD} \\ \mathbf{h} + \begin{pmatrix} 0 \\ 1 \end{pmatrix} & \text{otherwise} \end{cases}$
- $\mathbf{h}_0 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$

which counts the number of occurrences of positive and negative words.

To define *language models* using recurrent neural networks, we will use them as the encoder functions enc in our general language modeling framework (cf. §3.1). To connect Definition 5.1.1 with the general LM framework, we define the RNN encoding function.

Definition 5.1.2 (Recurrent Neural Encoding Function). Let $\mathcal{R} = (\Sigma, D, \mathbf{f}, \mathbf{h}_0)$ be a recurrent neural network. A **recurrent neural encoding function** $\text{enc}_{\mathcal{R}}$ is a representation function (cf. Definition 3.1.5) that recursively encodes strings of arbitrary lengths using its dynamics map \mathbf{f} :

$$\text{enc}_{\mathcal{R}}(\mathbf{y}_{<t+1}) \stackrel{\text{def}}{=} \mathbf{f}(\text{enc}_{\mathcal{R}}(\mathbf{y}_{<t}), y_t) \in \mathbb{R}^D \quad (5.4)$$

and

$$\text{enc}_{\mathcal{R}}(\mathbf{y}_{<1}) \stackrel{\text{def}}{=} \mathbf{h}_0 \in \mathbb{R}^D \quad (5.5)$$

Intuitively, an RNN \mathcal{R} takes an input string \mathbf{y} and encodes it with the encoding function $\text{enc}_{\mathcal{R}}$ by sequentially applying its dynamics map \mathbf{f} . The representations of individual prefixes (cf. Definition 2.3.6) of the input string are called hidden states.

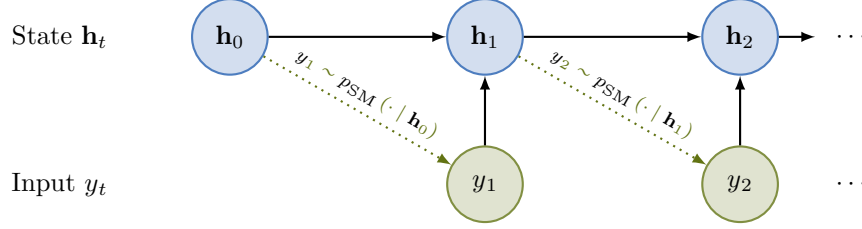
Definition 5.1.3 (Hidden State). Let $\mathcal{R} = (\Sigma, D, \mathbf{f}, \mathbf{h}_0)$ be an RNN. The hidden state $\mathbf{h}_t \in \mathbb{R}^D$ describes state of \mathcal{R} after reading y_t . It is recursively computed according to the dynamics map \mathbf{f} as follows:

$$\mathbf{h}_t \stackrel{\text{def}}{=} \text{enc}_{\mathcal{R}}(\mathbf{y}_{<t+1}) = \mathbf{f}(\mathbf{h}_{t-1}, y_t) \quad (5.6)$$

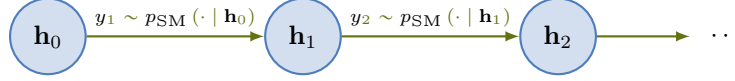
Example 5.1.5. The hidden states of the RNN from Example 5.1.2 are the individual values h_t , which, as t increases, approach $\sqrt{2}$.

Recurrent Neural Sequence Models

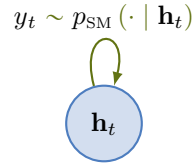
A recurrent neural network based on Definition 5.1.1 on its own does not yet define a sequence model, but simply a *context encoding function* $\text{enc}_{\mathcal{R}}: \bar{\Sigma}^* \rightarrow \mathbb{R}^D$. To define a sequence model based on an RNN, we simply plug in the RNN encoding function Definition 5.1.2 into the General language modeling framework from §3.1.



(a) An abstract depiction of how an RNN generates a string one symbol at a time. The hidden state \mathbf{h}_t summarizes the string $y_1 y_2 \dots y_t$ generated so far. The dotted lines denote the sampling steps.



(b) An abstract depiction of a generative RNN as an automaton.



(c) An abstract depiction of an RNN as a system updating the hidden state \mathbf{h}_t depending on the generated symbol y_t .

Figure 5.2: Different possible depictions of an abstract RNN model *generating* symbols. The way that the hidden states are updated based on the input symbol y_t is abstracted away.

Definition 5.1.4 (Recurrent neural sequence model). Let $\mathcal{R} = (\Sigma, D, \mathbf{f}, \mathbf{h}_0)$ be a recurrent neural network and $\mathbf{E} \in \mathbb{R}^{|\Sigma| \times D}$ a symbol representation matrix. A D -dimensional **recurrent neural sequence model** over an alphabet Σ is a tuple $(\Sigma, D, \mathbf{f}, \mathbf{E}, \mathbf{h}_0)$ defining the sequence model of the form

$$p_{SM}(y_t | \mathbf{y}_{<t}) \stackrel{\text{def}}{=} \mathbf{f}_{\Delta|\Sigma|-1}(\mathbf{E} \text{enc}_{\mathcal{R}}(\mathbf{y}_{<t}))_{y_t} = \mathbf{f}_{\Delta|\Sigma|-1}(\mathbf{E} \mathbf{h}_{t-1})_{y_t}. \quad (5.7)$$

By far the most common choice of the projection function is the softmax yielding the sequence model

$$p_{SM}(y_t | \mathbf{y}_{<t}) \stackrel{\text{def}}{=} \text{softmax}(\mathbf{E} \text{enc}_{\mathcal{R}}(\mathbf{y}_{<t}))_{y_t} = \text{softmax}(\mathbf{E} \mathbf{h}_{t-1})_{y_t}. \quad (5.8)$$

For conciseness, we will refer to RNN sequence models whose next-symbol probability distributions are computed using the softmax function as **softmax RNN sequence models**.

From this perspective, we see that RNNs are simply a special case of our general language modeling framework with parameterized representations of tokens $y \in \Sigma$ and the history $\mathbf{y} \in \Sigma^*$ (cf. §3.1)—an RNN simply defines how the encoding function enc is specified. The three figures from Fig. 5.1 are presented again with this probabilistic perspective in Fig. 5.2.

A Few More Definitions

In the following, we will often use the so-called one-hot encodings of symbols for concise notation. We define them here.

Definition 5.1.5 (One-hot encoding). Let Σ be an alphabet and $n : \Sigma \rightarrow \{1, \dots, |\Sigma|\}$ a bijection (i.e., an ordering of the alphabet, assigning an index to each symbol in Σ). A **one-hot encoding** $\llbracket \cdot \rrbracket$ is a representation function of the symbols in Σ which assigns the symbol $y \in \Sigma$ the $n(y)^{\text{th}}$ basis vector:

$$\llbracket y \rrbracket \stackrel{\text{def}}{=} \mathbf{d}_{n(y)}, \quad (5.9)$$

where here \mathbf{d}_n is the n^{th} canonical basis vector, i.e., a vector of zeros with a 1 at position n .

Example 5.1.6 (One-hot encoding). Let $\Sigma = \{\text{"large"}, \text{"language"}, \text{"models"}\}$ and $n = \{\text{"large"}: 1, \text{"language"}: 2, \text{"models"}: 3\}$. The one-hot encoding of the vocabulary is:

$$\llbracket \text{"large"} \rrbracket = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \llbracket \text{"language"} \rrbracket = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \llbracket \text{"models"} \rrbracket = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \quad (5.10)$$

Many specific variants of recurrent neural networks define the dynamics map \mathbf{f} in a specific way: the output of the function is some element-wise (non-linear) transformation of some “inner” function \mathbf{g} . The dynamics map of such an RNN is then the composition of \mathbf{g} and the non-linearity.

Definition 5.1.6 (Activation function). Let $\mathcal{R} = (\Sigma, D, \mathbf{f}, \mathbf{E}, \mathbf{h}_0)$ be an RNN. If the hidden states \mathbf{h}_t of the RNN are computed as

$$\mathbf{h}_t = \sigma(\mathbf{g}(\mathbf{h}_{t-1}, y)) \quad (5.11)$$

for some function $\mathbf{g}: \mathbb{R}^D \times \Sigma \rightarrow \mathbb{R}^D$ and some function $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ which is computed element-wise (that is, $\sigma(\mathbf{x})_d = \sigma(x_d)$ for all $d = 1, \dots, D$ and $\mathbf{x} \in \mathbb{R}^D$), we call σ an **activation function**.

This finishes our formal definition of recurrent neural networks. We next consider some of their theoretical properties, starting with tightness.

5.1.3 General Results on Tightness

We now discuss a general result on the tightness of recurrent neural sequence models, as defined in Definition 5.1.4. The analysis is straightforward and is a translation of the generic results on tightness (cf. §3.1.5) to the case of the norm of the hidden states of an RNN, \mathbf{h}_t as the encodings of the prefixes $\mathbf{y}_{\leq t}$, but it requires us to focus specifically on softmax RNN sequence models.

Theorem 5.1.1 (Tightness of Recurrent Neural Sequence Model). A softmax recurrent neural sequence model is tight if for all time steps t it holds that

$$s \|\mathbf{h}_t\|_2 \leq \log t, \quad (5.12)$$

where $s \stackrel{\text{def}}{=} \max_{y \in \Sigma} \|\mathbf{e}(y) - \mathbf{e}(\text{EOS})\|_2$.

Proof. This is simply a restatement of Theorem 3.1.6 for the case when enc takes the form of a general RNN encoding function, $\text{enc}_{\mathcal{R}}$. ■

Corollary 5.1.1 (RNNs with bounded dynamics maps are tight). A softmax recurrent neural sequence model $\mathcal{R} = (\Sigma, D, \mathbf{f}, \mathbf{h}_0)$ with a bounded dynamics map \mathbf{f} , i.e., with a dynamics map \mathbf{f} such that

$$|\mathbf{f}(\mathbf{x})_d| \leq M \quad (5.13)$$

for some $M \in \mathbb{R}$, for all $d = 1, \dots, D$ and all $\mathbf{x} \in \mathbb{R}^D$, is tight.

Proof. If the dynamics map is bounded, the norm of the hidden state, $\|\mathbf{h}_t\|_2$, is bounded as well. This means that the left-hand-side of Eq. (5.12) is constant with respect to t and the condition holds trivially. ■

A special case of Corollary 5.1.1 is RNNs with bounded activation functions (cf. Definition 5.1.6). Those are tight if the activation function itself is bounded. This implies that all standard sigmoid and tanh activated recurrent neural networks are tight. However, the same does not hold for RNNs with unbounded activation functions, which have lately been more popular (one of the reasons for this is the vanishing gradient problem (Glorot et al., 2011)).

Example 5.1.7 (RNNs with unbounded activation functions may not be tight). A very popular unbounded activation function is the so-called **rectified linear unit** (ReLU), defined as

$$\text{ReLU}(x) \stackrel{\text{def}}{=} \max(0, x). \quad (5.14)$$

This function is clearly unbounded.

Now suppose we had the following RNN over the simple alphabet $\Sigma = \{a\}$.

$$\mathbf{h}_t = \mathbf{h}_{t-1} + (1), \quad (5.15)$$

initial state

$$\mathbf{h}_0 = (0) \quad (5.16)$$

and the output matrix

$$\mathbf{E} = \begin{pmatrix} -1 \\ 1 \end{pmatrix} \quad (5.17)$$

where the top row of \mathbf{E} computes the logit of the EOS symbol and the bottom one the one of a . It is easy to see that

$$\mathbf{h}_t = (t). \quad (5.18)$$

This already does not look promising for tightness—the norm of the hidden state, which is, in this case, $\|\mathbf{h}_t\| = t$, and is, therefore, increasing at a much higher rate than $\mathcal{O}(\log t)$ required by Theorem 5.1.1. We encounter a similar hint against tightness if we compute the conditional probabilities of the EOS symbol and the symbol a .

$$p_{SM}(\text{EOS} \mid \mathbf{y}_{<t}) = \text{softmax} \left(\begin{pmatrix} -1 \\ 1 \end{pmatrix} (t-1) \right)_{\text{EOS}} = \frac{\exp[-t+1]}{\exp[-t+1] + \exp[t-1]} \quad (5.19)$$

$$p_{SM}(\text{EOS} \mid \mathbf{y}_{<t}) = \text{softmax} \left(\begin{pmatrix} -1 \\ 1 \end{pmatrix} (t-1) \right)_a = \frac{\exp[t-1]}{\exp[-t+1] + \exp[t-1]} \quad (5.20)$$

The probability of ending the string at time step t is, therefore

$$p_{SM}(\text{EOS} \mid \mathbf{y}_{<t}) = \frac{1}{1 + \exp[2(t-1)]}. \quad (5.21)$$

Intuitively, this means that the probability of ending the string (generating EOS) diminishes rapidly with t —in this case much faster than any diverging sum required by Theorem 2.5.3. All signs, thus, point towards the RNN from Eq. (5.15) not being tight. Indeed, for this specific case, one can show using some algebraic manipulations that

$$\sum_{\mathbf{y} \in \Sigma^*} p_{LN}(\mathbf{y}) = \sum_{n \in \mathbb{N}_{\geq 0}} p_{LN}(a^n) < 0.15 \quad (5.22)$$

where p_{LN} is the locally normalized model induced by the RNN. This means that the RNN from Eq. (5.15) assigns less than 0.15 probability to finite strings—all other probability mass leaks to infinite sequences.

Example 5.1.8 (RNNs with unbounded activation functions can still be tight). Example 5.1.7 showed that RNNs with unbounded activation functions can indeed result in non-tight sequence models. However, this is not necessarily the case, as this simple modification of the RNN from Example 5.1.7 shows. The only aspect of the RNN that we modify is the output matrix \mathbf{E} , which we change by flipping its rows:

$$\mathbf{E} = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad (5.23)$$

Now the probability of ending the string at time step t is

$$p_{SM}(\text{EOS} \mid \mathbf{y}_{<t}) = \frac{\exp[t-1]}{\exp[-t+1] + \exp[t-1]} = \frac{1}{\exp[-2(t-1)] + 1}. \quad (5.24)$$

Compared to Eq. (5.21), the probability of EOS in Eq. (5.24) does not diminish. Indeed, since $\frac{1}{\exp[-2(t-1)] + 1} > \frac{1}{2}$ for all t , the sum

$$\sum_{t=0}^{\infty} p_{SM}(\text{EOS} \mid \mathbf{y}_{<t}) \quad (5.25)$$

diverges, which, according to Proposition 2.5.6 implies that the sequence model is tight.

5.1.4 Elman and Jordan Networks

The characterization of dynamics maps we gave in Definition 5.1.4 allows for \mathbf{f} to be an arbitrary mapping from the previous state and the current input symbol to the new state. In this section, we introduce two seminal and particularly simple parameterizations of this map—the simplest recurrent neural sequence models. We term them Elman sequence models and Jordan sequence models, as each is inspired by architectures proposed by Elman (1990) and Jordan (1986), respectively. The definitions we present here are slightly different than those found in the original works—most notably, both Elman and Jordan networks were originally defined for *transduction* (mapping an input string to an output string, as with translation) rather than language modeling.

Put simply, these two models *restrict* the form of the dynamics map \mathbf{f} in the definition of an RNN (cf. Definition 5.1.1). They define particularly simple relationships between the subsequent hidden states, which are composed of affine transformations of the previous hidden state and the representation of the current input symbol passed through a non-linear activation function (cf. Definition 5.1.6). The affine transformations are performed by different matrices and bias vectors—the parameters of the model (cf. Assumption 3.2.2)—each transforming a separate part of the input to the dynamics map.

Definition 5.1.7 (Elman Sequence Model (Elman, 1990)). *An **Elman sequence model** $\mathcal{R} = (\Sigma, D, \mathbf{U}, \mathbf{V}, \mathbf{E}, \mathbf{b}_h, \mathbf{h}_0)$ is a D -dimensional recurrent neural sequence model over an alphabet Σ with the following dynamics map*

$$\mathbf{h}_t = \sigma(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{V}\mathbf{e}'(y_t) + \mathbf{b}_h). \quad (5.26)$$

Here, $\mathbf{e}' : \Sigma \rightarrow \mathbb{R}^R$ is the input symbol **embedding function** which represents each symbol $y \in \Sigma$ as a R -dimensional vector and σ is an element-wise non-linearity.⁶ $\mathbf{b}_h \in \mathbb{R}^D$, $\mathbf{U} \in \mathbb{R}^{D \times D}$, and $\mathbf{V} \in \mathbb{R}^{D \times R}$.

Due to its simplicity, the Elman RNN is also known as the *vanilla RNN* variant, emphasizing it is one of the most fundamental variants of the framework.

On the symbol representations. Notice that in Eq. (5.26), the input symbols y_t are first transformed into their vector representations $\mathbf{e}'(y_t)$ and then additionally linearly transformed using the matrix \mathbf{V} . This results in an over-parametrized network—since the symbols are already embedded using the representation function \mathbf{e}' , the matrix \mathbf{V} is theoretically superfluous and could be replaced by the identity matrix. However, the matrix \mathbf{V} could still be useful if the representations $\mathbf{e}'(y_t)$ are fixed—in this case, the matrix can be used by the RNN to transform the representations during training to fit the training data better. This is especially useful if the symbol representations $\mathbf{e}'(y_t)$ already represent the input symbols in a compact representation space in which the parameters can be shared across different symbols. Alternatively, we could represent the symbols using their one-hot encodings, i.e., $\mathbf{e}'(y_t) = \llbracket y_t \rrbracket$, in which case the columns of the matrix \mathbf{V} would correspond to the symbol representations (analogously to the representation matrix \mathbf{E} from Eq. (3.46)). However, notice that in this case, the representations on the symbols do not share any parameters, and each column of the matrix is therefore an unconstrained vector. Such matrix-lookup-based *input* symbol representations from Eq. (5.26) are sometimes **tied**, i.e., $\mathbf{e}' = \mathbf{e}(\cdot)$, with the *output* symbol representations from the embedding matrix \mathbf{E} in the definition of the sequence model induced by an RNN (cf. Definition 3.1.11 and Eq. (5.7)).

However, embedding tying is non-essential to representation-based LMs. The input symbol embedding function can always be chosen independently with the output symbol embedding function Definition 3.1.6.

The Jordan network is somewhat different in that it feeds the *output* logits computed through the output matrix \mathbf{E} into the computation of the next state, and not directly the hidden state.

Definition 5.1.8 (Jordan Sequence Model (Jordan, 1986)). *A **Jordan sequence model** is a D -dimensional recurrent neural sequence model over an alphabet Σ with the following dynamics map*

$$\mathbf{h}_t = \sigma(\mathbf{U}\mathbf{r}_{t-1} + \mathbf{V}\mathbf{e}'(y_t) + \mathbf{b}_h) \quad (5.27)$$

$$\mathbf{r}_t = \sigma_o(\mathbf{E}\mathbf{h}_t) \quad (5.28)$$

⁶The symbol representations $\mathbf{e}'y$ are often also referred to as **static symbol embeddings** because they do not depend on the string surrounding or preceding y . Here, we treat them as any other parameters of the model which can be learned using gradient-based learning (cf. §3.2). However, note that learning good static embeddings was a very active field before the emergence of large end-to-end systems we see today. Very popular examples include Word2Vec (Mikolov et al., 2013), GloVe (Pennington et al., 2014), and FastText (Bojanowski et al., 2017).

Again, $\mathbf{e}' : \Sigma \rightarrow \mathbb{R}^R$ is the input symbol embedding function which represents each symbol $y \in \Sigma$ as a R -dimensional vector while σ and σ_o are element-wise non-linearities. $\mathbf{b}_h \in \mathbb{R}^D$, $\mathbf{U} \in \mathbb{R}^{D \times D}$, and $\mathbf{V} \in \mathbb{R}^{D \times R}$.

Notice that the hidden state \mathbf{h}_t in Eq. (5.27) is not computed based on the previous hidden state \mathbf{h}_{t-1} , but rather on the transformed outputs \mathbf{r}_{t-1} —this is analogous to feeding back in the logits computed in Eq. (5.7) into the computation of \mathbf{h} rather than the previous hidden state. The sequence model induced by a Jordan network is then directly induced by the logits \mathbf{r}_t (i.e., the conditional probabilities are computed by putting \mathbf{v}_t through the softmax).

In both architectures, the activation function σ can be any suitable element-wise function. The canonical choices for it have been the sigmoid and tanh functions, however, a more common choice nowadays is the ReLU function or any of its more modern variants.⁷

Since we will refer to the individual matrices defining the dynamics maps in Elman and Jordan networks quite a lot in the next subsections, we give them specific names. The matrix \mathbf{U} , which linearly transforms the previous hidden state (or the output) is the **recurrence matrix**. The matrix \mathbf{V} , which linearly transforms the representations of the input symbol, is called the **input matrix**. Lastly, the matrix which linearly transforms the hidden state before computing the output values \mathbf{r}_t with an activation function is called the **output matrix**. \mathbf{b}_h is the hidden **bias vector**.

Tightness of Elman and Jordan Recurrent Neural Networks As a simple corollary of Corollary 5.1.1, we can characterize the tightness of Elman and Jordan recurrent neural networks as follows.

Corollary 5.1.2 (Tightness of simple RNNs). *Elman and Jordan RNNs with a bounded activation function σ and the softmax projection function are tight.*

⁷See Goodfellow et al. (2016, §6.3.1) for an overview of modern activation functions used in neural networks.

5.1.5 The Vanishing and Exploding Gradient Problem

Let us use the following examples to motivate this subsection.

Example 5.1.9 (Long-term dependency in natural language). *Considering one of the sentences from Example 4.2.1.*

“The cat the dog the mouse the rat frightened startled barked at likes to cuddle.”

Questions: Who likes to cuddle? Who barked? What did the mouse do?

Example 5.1.10 (Long-term dependency in formal language). *Consider the language $L = \{ab^n a \mid n \in \mathbb{N}\} \cup \{cb^n c \mid n \in \mathbb{N}\}$. Does the string “ $ab^n c$ ” belong to the language?*

To be able to answer these questions, a language model must capture **long-term dependencies**—the relationships between words that are separated by many other words. How to model such dependencies has long been a key issue for language models. We have already shown that RNNs have the ability to model inputs of unbounded length. Yet, this ability should not be taken for granted. It is to be learned from data with considerable effort.

In most cases, we have specifically chosen the dynamics map \mathbf{f} so that the function is *differentiable*. Because following the gradient of the model’s log-likelihood with respect to the parameters is a useful search technique for finding a good model, it is worthwhile to understand how the gradient behaves. Unfortunately, the gradient of an Elman network is a particularly unstable quantity. As a function of the length of a training exemplar, in most cases, it tends to either zero or infinity. This fact, known as the **vanishing and exploding gradient problem**, respectively, is the motivation for many improvements to the Elman network.

In the following discussion, we analyze the gradient behavior of the parameters of an RNN using a generic differentiable $\ell: \Theta \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ as the training objective, denoting θ as some learnable parameters in the RNN. Vanishing/exploding gradient refers to the cases in which $\frac{\partial \ell}{\partial \theta}$ is consistently close to $\mathbf{0}$ or extremely large in scale, regardless of the input training exemplar. It is important to keep in mind that vanishing/exploding gradient is different with the extreme gradients on parameters when the model receives cliché or unexpected input—the latter one is desirable for training while the former one is not.

Before we start, a few assumptions must be made for our analysis to be meaningful.

Fact 5.1.1 (RNN loss function). *The loss function ℓ of an RNN can be written as a function of $\mathbf{h}_0, \mathbf{h}_1, \dots, \mathbf{h}_T, \mathbf{E}$ —the sufficient statistics of an input sequence \mathbf{y} .*

$$\ell(\theta, \theta^*) \stackrel{\text{def}}{=} \ell(\mathbf{h}_0, \mathbf{h}_1, \dots, \mathbf{h}_T, \mathbf{E})$$

Therefore, its gradient with regard to the parameters θ can be decomposed using the following application of the chain rule:

$$\frac{\partial \ell}{\partial \theta} = \sum_{t=0}^T \frac{\partial \ell}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \theta} + \frac{\partial \ell}{\partial \mathbf{E}} \frac{\partial \mathbf{E}}{\partial \theta} \quad (5.29)$$

The reasoning behind this fact is that the hidden states $\mathbf{h}_0, \mathbf{h}_1, \dots, \mathbf{h}_T, \mathbf{E}$ are *sufficient* for modeling any desired output given \mathbf{y} —they have taken all of the input string and the learned parameters of the model into consideration. In the context of language modeling, this fact can be rephrased into “ \mathbf{h}_{t-1} and \mathbf{E} are the sufficient statistics for $p_{SM}(y_t \mid \mathbf{y}_{<t}) = p_{\theta}(y_t \mid \mathbf{y}_{<t})$ —the likelihood of the symbol y_t ” (cf. Definition 5.1.4, the softmax for modeling the distribution of y_t only takes \mathbf{h}_{t-1} and \mathbf{E} as input⁸).

In Eq. (5.29), the term $\frac{\partial \ell}{\partial \mathbf{h}_t}$ is a gradient of a one-time application of the loss function, which is stable and easy to evaluate. Then, we can narrow the scene of vanishing/exploding gradient down to the computation of $\frac{\partial \mathbf{h}_t}{\partial \theta}$.

⁸This conclusion follows by Neyman-Fisher factorization theorem (Halmos and Savage, 1949), since the conditional likelihood can be represented as $p_{\theta}(y_t \mid \mathbf{y}_{<t}) = g_{\theta}(\mathbf{h}_{t-1}, \mathbf{E})$.

Without loss of generality, let us look at the case when $t = T$. $\frac{\partial \mathbf{h}_T}{\partial \theta}$ can be further expanded into

$$\begin{aligned}
 \frac{\partial \mathbf{h}_T}{\partial \theta} &= \frac{\partial \mathbf{f}(\mathbf{h}_{T-1}, y_T)}{\partial \theta} \\
 &= \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \frac{\partial \mathbf{h}_{T-1}}{\partial \theta} + \frac{\partial_{\bar{\mathbf{h}}} \mathbf{h}_T}{\partial \theta} \\
 &= \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \left(\frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{h}_{T-2}} \frac{\partial \mathbf{h}_{T-2}}{\partial \theta} + \frac{\partial_{\bar{\mathbf{h}}} \mathbf{h}_{T-1}}{\partial \theta} \right) + \frac{\partial_{\bar{\mathbf{h}}} \mathbf{h}_T}{\partial \theta} \\
 &= \dots \\
 &= \sum_{t=0}^T \left(\prod_{i=t+1}^T \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} \right) \frac{\partial_{\bar{\mathbf{h}}} \mathbf{h}_t}{\partial \theta} \\
 &= \sum_{t=0}^T \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} \frac{\partial_{\bar{\mathbf{h}}} \mathbf{h}_t}{\partial \theta}
 \end{aligned}$$

$\frac{\partial_{\bar{\mathbf{h}}} \mathbf{h}_T}{\partial \theta}$ refers to the partial derivative of $\mathbf{h}_T = \mathbf{f}(\mathbf{h}_{T-1}, y_T)$ w.r.t. θ , while treating \mathbf{h}_{T-1} as a constant. I.e., the gradient does not backpropagate into \mathbf{h}_{T-1} . Just like $\frac{\partial \ell}{\partial \mathbf{h}_t}$ in Eq. (5.29), $\frac{\partial_{\bar{\mathbf{h}}} \mathbf{h}_T}{\partial \theta}$ is the derivative of an one-time application of a function, which is easy to evaluate.

On the contrary, the deeply nested $\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t}$ can be very unstable. Mathematically speaking, an RNN is little more than a composition function. It functions by repeatedly applying the same dynamics map on the hidden state at each time step.

Consider T applications of an RNN's dynamics maps

$$\mathbf{h}_T = \mathbf{f}(\dots \mathbf{f}(\mathbf{h}_0, y_1) \dots, y_T)$$

Recall Example 5.1.10, when deciding whether " $ab^n c$ " belongs to $L = \{ab^n a \mid n \in \mathbb{N}\} \cup \{cb^n c \mid n \in \mathbb{N}\}$ or not, the equation above is instantiated to

$$\mathbf{h}_T = \mathbf{f}(\underbrace{\dots \mathbf{f}(\mathbf{h}_0, a)}_{\mathbf{h}_1} \dots, c)$$

The decision can only be made after \mathcal{R} reads the last symbol $y_T = c$, i.e., when computing \mathbf{h}_T . And most importantly, the computation of \mathbf{h}_T should take the initial symbol y_1 into account— $y_1 = a$ or $y_1 = c$ matter a lot.

We use the sufficient statistics \mathbf{h}_1 as the proxy of y_1 . And examine how much \mathbf{h}_T is affected when \mathbf{h}_1 changes. The computation of \mathbf{h}_T can be regarded as a vector-valued function, in which \mathbf{h}_1 is a variable. We are interested in the Jacobian $\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_1}$, which is the best linear approximation of the change in \mathbf{h}_T when giving a small perturbation to \mathbf{h}_1 . Mathematically speaking, we have

$$\Delta \mathbf{h}_T \approx \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_1} \Delta \mathbf{h}_1$$

With the Jacobian matrix, we can now *quantitatively* analyze long-term dependency.

Again, without loss of generality, we analyze $\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0}$ to get a sense of the gradient effects of long-distance dependencies between hidden states.

Now that we have located the Vanishing/exploding gradient arises when the L^p -norm $\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_p$ being almost 0 or extremely large—when it's 0, long-distance dependencies are not reflected in the gradient updates; when it's extremely large, long-distance dependencies are overemphasized in at least one dimension. One term in Eq. (5.29), the gradient on some parameter θ :

$$\frac{\partial \ell}{\partial \mathbf{h}_T} \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \frac{\partial \mathbf{h}_0}{\partial \theta}$$

could be extremely unstable to evaluate. Because an overflow is likely to happen while computing the matrix-vector product.

We aim to bound the L^1 -norm $\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1$ from above and below in the following discussions. The L^1 -norm is defined as follows

Definition 5.1.9 (Matrix L^1 -norm). Suppose $\mathbf{A} \in \mathbb{R}^{D \times D}$ with column vectors $\{\mathbf{v}_i\}_{i=1,2,\dots,D}$, the L^1 -norm of \mathbf{A} is defined as

$$\|\mathbf{A}\|_1 = \max_{i \in \{1,2,\dots,D\}} \|\mathbf{v}_i\|_1 \quad (5.30)$$

In theory, any norm could work for our discussion. But L^1 -norm makes it easier for us to derive the results for the lower bound (Lemma 5.1.1). For vanishing or exploding gradient, we will respectively give the sufficient condition for the upper bound to *decrease* exponentially or the lower bound to *grow* exponentially.

Theorem 5.1.2 (Vanishing Gradient). Suppose the L^1 -norm of the Jacobian $\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\|_1 \leq C < 1$ is bounded above by a constant, itself less than 1. Then $\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1$ decreases to $\mathbf{0}$ exponentially.

Proof. Theorem 5.1.2 is a direct consequence of the matrix norm submultiplicativity:

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1 = \left\| \prod_{t=1}^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\|_1 \leq \prod_{t=1}^T \left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\|_1 \quad (5.31)$$

Suppose $\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\|_1 \leq C < 1$, we have

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1 \leq \prod_{t=1}^T \left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\|_1 \leq C^T \quad (5.32)$$

which decreases to $\mathbf{0}$ exponentially. ■

Now, let's turn to another extremity of gradient behaviors.

Theorem 5.1.3 (Exploding Gradient). Suppose the determinant of the Jacobian $\det\left(\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}\right) \geq C > 1$ is bounded below by a constant, itself greater than 1. Then $\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1$ increases to ∞ exponentially.

The proof of Theorem 5.1.3 is a bit trickier than the one of Theorem 5.1.2. To begin with, there is no inequality at hand for us to give a lower bound to the norm of a matrix product.

However, recall the multiplicativity of determinant: $\det(\mathbf{AB}) = \det(\mathbf{A})\det(\mathbf{B})$ for square matrices \mathbf{A}, \mathbf{B} . It hints that the determinant can be used as a bridge in lower bounding the product of $\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}$. Meanwhile, the determinant itself can serve as a lower bound for the L^1 -norm by the following inequality:

Fact 5.1.2 (Hadamard's Inequality). The following inequality holds for $\mathbf{A} \in \mathbb{R}^{D \times D}$ with column vectors $\{\mathbf{v}_i\}_{i=1,2,\dots,D}$

$$\det(\mathbf{A}) \leq \prod_{i=1}^D \|\mathbf{v}_i\|_2 \quad (5.33)$$

The equality is achieved if and only if \mathbf{v}_i are orthogonal. For detailed proof, we refer the readers to [Horn and Johnson \(2012, Corollary 7.8.3\)](#).

By Fact 5.1.2, we can bound the L^1 -norm of a matrix \mathbf{A} from below using its determinant $\det(\mathbf{A})$.

Lemma 5.1.1 (Determinantal lower bound on matrix L^1 norm). Let $\mathbf{A} \in \mathbb{R}^{D \times D}$, we have

$$\det(\mathbf{A}) \leq \|\mathbf{A}\|_1^D$$

Proof. Let $\mathbf{A} \in \mathbb{R}^{D \times D}$, we have

$$\begin{aligned} \det(\mathbf{A}) &\leq \prod_{i=1}^D \|\mathbf{v}_i\|_2 && \text{(Hadamard's inequality)} \\ &\leq \prod_{i=1}^D \|\mathbf{v}_i\|_1 && \text{(Cauchy-Schwarz inequality)} \\ &\leq \|\mathbf{A}\|_1^D && \text{(Property of } L^1 \text{ norm)} \end{aligned}$$

The first equality is achieved if and only if \mathbf{v}_i are orthogonal. The second equality is achieved if and only if \mathbf{v}_i are $\mathbf{0}$ or parallel with standard basis vectors. The third equality is achieved if and only if all $\|\mathbf{v}_i\|_1$ are the same. ■

Note that when $\det(\mathbf{A}) \leq 0$, Lemma 5.1.1 becomes a trivial inequality $\|\mathbf{A}\|_1 \geq 0$.

With Lemma 5.1.1, we can now present the proof for Theorem 5.1.3:

Proof of Theorem 5.1.3. By Lemma 5.1.1, we give a lower bound on $\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1$:

$$\begin{aligned} \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1 &\geq \left(\det \left(\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right) \right)^{\frac{1}{D}} \\ &= \left(\det \left(\prod_{t=1}^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right) \right)^{\frac{1}{D}} \\ &= \prod_{t=1}^T \left(\det \left(\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right) \right)^{\frac{1}{D}} \end{aligned} \quad (5.34)$$

Suppose $\det \left(\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right) \geq C > 1$, we have

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1 \geq C^{\frac{T}{D}}$$

which increases to ∞ exponentially as $T \rightarrow \infty$. ■

Combining the upper bound (Eq. (5.31)) and the lower bound (Eq. (5.34)), we have:

Corollary 5.1.3. $\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1$ satisfies the following inequalities:

$$\prod_{t=1}^T \left(\det \left(\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right) \right)^{\frac{1}{D}} \leq \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1 \leq \prod_{t=1}^T \left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\|_1 \quad (5.35)$$

Up to this point, we have concluded our analysis of the gradient of RNNs. A hot takeaway is that gradient-based optimization can be extremely unstable for long-term dependencies. We must carefully choose the dynamics map \mathbf{f} so as to avoid the vanishing/exploding gradient conditions.

To understand how Theorems 5.1.2 and 5.1.3 matter in practice, let us take the Elman sequence model as an example to instantiate the above bounds. Recall its dynamics map \mathbf{f} :

$$\begin{aligned} \mathbf{h}_0 &= \mathbf{0} \\ \mathbf{h}_t &= \mathbf{f}(\mathbf{h}_{t-1}, y_{t-1}) = \sigma(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{V}\mathbf{e}'(y_{t-1}) + \mathbf{b}_h) \end{aligned} \quad t \geq 1 \quad (5.36)$$

where σ is a non-linearity applied to the input element-wise. The step-wise gradient for one step of hidden state update is:

$$\begin{aligned} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} &= \frac{\partial \mathbf{f}(\mathbf{h}_{t-1}, y_{t-1})}{\partial \mathbf{h}_{t-1}} \\ &= \text{diag}(\sigma'(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{V}\mathbf{e}'(y_{t-1}) + \mathbf{b}_h))\mathbf{U} \\ &\stackrel{\text{def}}{=} \text{diag}(\sigma'_t)\mathbf{U} \end{aligned} \quad (5.37)$$

Fact 5.1.3. For Elman networks, $\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1$ satisfies the following inequalities

$$\prod_{t=1}^T (\det(\text{diag}(\sigma'_t)\mathbf{U}))^{\frac{1}{D}} \leq \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1 \leq \prod_{t=1}^T \|\text{diag}(\sigma'_t)\mathbf{U}\|_1 \quad (5.38)$$

Vanishing Gradient

To find out when vanishing gradient could happen, we aim to find the $\text{diag}(\sigma'_t)\mathbf{U}$ that makes the upper bound decrease exponentially with the length of the input.

The derivative of activation function σ' is bounded in many cases (e.g., ReLU, sigmoid). Denoting M as the upper bound such that $\forall x \in \mathbb{R}, |\sigma'(x)| \leq M$, we have

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1 \leq (M \|\mathbf{U}\|_1)^T \quad (5.39)$$

The term $(M \|\mathbf{U}\|_1)^T$ exponentially shrinks to 0 when $M \|\mathbf{U}\|_1 < 1$. The hidden states corresponding to the earlier part of the input sequence will receive very limited gradient updates.

Exploding Gradient

A weaker result can immediately be derived by inverting the proof for the upper bound: $M \|\mathbf{U}\|_1 \geq 1$ is the *necessary* condition for exploding gradient. Since otherwise, the gradient would vanish.

However, we are not satisfied with this conclusion. By Fact 5.1.3, we present a sufficient condition that leads to exploding gradient⁹.

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1 \geq \left(\prod_{t=1}^T \det(\text{diag}(\sigma'_t)) \det(\mathbf{U}) \right)^{\frac{1}{D}} \quad (5.40)$$

When $\frac{T}{D}$ is sufficiently large, and $\det(\mathbf{U}) \det(\text{diag}(\sigma'_t)) > 1$, an exploding gradient will appear.

For instance, let the activation function $\sigma = \text{ReLU}$, $\mathbf{U}\mathbf{h}_t + \mathbf{V}\mathbf{e}'(y_t) + \mathbf{b}_h > 0$ for all $t \geq 1$, and $\det(\mathbf{U}) = M > 1$. Then, $\det(\text{diag}(\sigma'_t)) = 1$, for all $t \in \{1, 2, \dots, T\}$. We have

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|_1 \geq M^{\frac{T}{D}} \quad (5.41)$$

which exponentially increases to ∞ as $T \rightarrow \infty$.

Why is vanishing/exploding gradient a problem?

The computing frameworks we employ to implement a language model often have finite precision. The parameters and intermediate variables during computing are represented by floating-point numbers, among which `float32` is most commonly used.

The minimal value above zero that a `float32` number can represent is 1.2×10^{-38} , while the maximal is 3.4×10^{38} . All values outside this range will be truncated to one of the three values: 0, NaN, or inf.

When vanishing gradient appears, gradient signals from long-term dependencies are extremely small in scale. Gradient signals could also be rounded off to zero due to the finite float precision. When making decisions, the trained model tends to lose track of long-distance dependencies and only relies on the nearest neighbors of the hidden states \mathbf{h}_t it chooses to use. On the other hand, a gradient that is too large in scale could lead to NaN updates during training, interrupting the training process.

In §5.1.6, we will discuss some techniques for countering vanishing/exploding gradient.

5.1.6 Variations on Recurrent Networks

In the previous sections we introduced the two simplest RNN variants: the Elman (Elman, 1990) and Jordan (Jordan, 1997) networks. Even though such simple RNNs in theory are all we need to model any computable language, empirically those architectures face many challenges. One of the biggest are the vanishing and exploding gradient problems (Hochreiter and Schmidhuber, 1997), which in practice is linked with the issue of learning long-term dependencies in language.

⁹Our analysis does not make assumptions on the choice of non-linearity, which makes itself more general than the analysis of Pascanu et al. (2013) that assumes identity activation.

In this subsection, we expand our repertoire of RNN variants by going beyond the simple recurrent dynamics defined by the Elman and Jordan update rules. To do so, we take a step back and return to Definition 5.1.1 of a recurrent neural network \mathcal{R} as the tuple $(\Sigma, D, \mathbf{f}, \mathbf{h}_0)$. We will define more elaborate dynamics maps \mathbf{f} which both aim to tackle some of the (empirically encountered) challenges of simpler variants as well as improve some theoretical aspects of the networks. Importantly, keep in mind that the only aspect of the RNN we will strive to modify is the dynamics map—that is, the mapping from \mathbf{h}_{t-1} to \mathbf{h}_t . Given a hidden state, the definition of a sequence model will remain identical.

A common component of the more complex dynamics maps we explore in this section is the gating mechanism, which is why we start with it.

Gating

The update equations of Elman and Jordan RNNs define relatively simple transformations of the hidden states as an affine transformation of the previous hidden state and the new input, followed by some form of non-linearity. In this sense, the interaction between the previous hidden state and the input symbol is relatively limited—the hidden state is transformed by the recurrence matrix \mathbf{U} at every time step invariant to the input symbol being read. To see why this could be a limiting factor, consider the following example.

Example 5.1.11. Consider the language $L = \{a^n b^n c^n x a^m b^m c^m \mid n, m \in \mathbb{N}_{\geq 0}\}$. It intuitively consists of two-part strings, where the two parts are separated by a symbol x . The part on the left side of x contains a sequence of n a 's followed by n b 's, which is followed by n c 's. The substring on the right side of x contains a sequence of m a 's which is again followed by m b 's, and later by m c 's. Both parts of the string can be arbitrarily long, and, intuitively, to correctly recognize a string in this language, a computational model has to keep the information about the number of a 's while reading in b 's to be able to ensure there is a correct number of c 's as well. This creates a long-term dependency across the entire block of b 's. However, notice that, after reading the symbol x , the information about the number of a 's becomes irrelevant to the recognition of the string: the model can, therefore, discard it and solely focus on modeling the rest of the string, which again requires keeping track of the number of the symbol occurrences. In other words: after a certain amount of time, previous information becomes irrelevant, and we may want to design a network that is able to select which information is important to keep around.

To enable richer interaction between the transformation of the RNN hidden state and the input symbols, we introduce the gating mechanism. Intuitively, the gating mechanism enables more fine-grained control over the transformations of the hidden state by “selecting” which aspects of the hidden state should be retained, which should be modified, and which should be deleted—in general, based on both the previous hidden state as well as the current input symbol. Such transformations are defined using gates and gating functions.

Definition 5.1.10 (Gate). A **gate** is a real-valued vector $\mathbf{g} \in \mathbb{R}^D$, such that $g_d \in [0, 1]$ for all $d \in \{1, \dots, D\}$. Gates are computed using **gating functions**, i.e., functions whose outputs live in $[0, 1]^D$.

The fact that every dimension in a gate \mathbf{g}_t takes a value between 0 and 1 invites a natural interpretation of the values as *soft switches*, analogously to how switches are used in the electrical engineering context. Intuitively, in the context of RNNs, where the information is passed around in the hidden states \mathbf{h}_t , a gate of the same dimensionality as the hidden state can control which aspects (dimensions) of the hidden state should be forgotten (switched off) or and which ones retained (kept on)—a gate value close to 0 can be interpreted as a signal that the information captured in the corresponding dimension of the hidden state should be “forgotten”, and a gate value close to 1 as the opposite. Such modifications of the hidden state can be performed using an *element-wise multiplication* of the hidden state \mathbf{h} and the gate \mathbf{g} , which we denote with $\mathbf{h} \odot \mathbf{g}$.

Importantly, the gates can be computed based on the information about the string seen so far as well as the new input symbol—this means that the decision on what should be remembered and what should be forgotten can be made for each situation individually. This allows RNN variants using gating to implement mechanisms to tackle challenges as the one described in Example 5.1.11. Furthermore this not only enables RNNs to selectively keep information about the string, but also combat the vanishing and exploding gradient problems (Hochreiter and Schmidhuber, 1997, Appendix 2). We next consider two of the best-known gated RNNs: Long Short-Term Memory and Gated Recurrent Unit networks.

Long Short-term Memory Networks

Long Short-term Memory Networks (LSTM, Hochreiter and Schmidhuber, 1997) are perhaps the best-known type of a gated recurrent network. They were introduced specifically to combat the vanishing gradient problem in the famous paper with more than 80 000 citations. The somewhat unusual name comes from connections to human memory, in which short-term memory is characterized by evanescent neural activations, and long-term memory is based on the growth and structural change in neuron connections (Hebb, 1949).

The LSTM unit. LSTM RNNs are built from the LSTM units, which implement the RNN dynamics map and therefore perform the RNN update step. To transform the hidden state \mathbf{h}_t at each time step, an LSTM network additionally keeps *another* running summary of the string $\mathbf{y}_{<t}$, on which the recurrent update depends—this is the so-called memory cell which we will denote by \mathbf{c}_t . Informally, one can think of the context as the information needed to decide on *how* to transform the hidden state at each individual time step, depending on the input string. The formal definition of the LSTM cell is the following.

Definition 5.1.11 (Long Short-Term Memory). *A **long short-term memory unit** is a recurrent neural network with the dynamics map defined through the following sequence of computations:*

$$\begin{aligned}
 \mathbf{i}_t &= \sigma(\mathbf{U}^i \mathbf{h}_{t-1} + \mathbf{V}^i \mathbf{e}' y_t + \mathbf{b}^i) && \text{(input gate)} \\
 \mathbf{f}_t &= \sigma(\mathbf{U}^f \mathbf{h}_{t-1} + \mathbf{V}^f \mathbf{e}' y_t + \mathbf{b}^f) && \text{(forget gate)} \\
 \mathbf{o}_t &= \sigma(\mathbf{U}^o \mathbf{h}_{t-1} + \mathbf{V}^o \mathbf{e}' y_t + \mathbf{b}^o) && \text{(output gate)} \\
 \mathbf{g}_t &= \tanh(\mathbf{U}^g \mathbf{h}_{t-1} + \mathbf{V}^g \mathbf{e}' y_t + \mathbf{b}^g) && \text{(candidate vector)} \\
 \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t && \text{(memory cell)} \\
 \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t) && \text{(hidden state)}
 \end{aligned}$$

$\mathbf{i}_t, \mathbf{f}_t, \mathbf{o}_t$ are the **input**, **forget**, and **output** gates, \mathbf{c}_t is the **memory cell** vector, and \mathbf{g}_t is the **candidate vector**. Here, σ refers to the original sigmoid function.

As we can see, the update rule of an LSTM network is considerably more complex than that of an Elman RNN. It is also computationally more expensive, as it involves more matrix multiplications. However, LSTMs have consistently shown improved performance compared to vanilla RNNs and are therefore considered together with GRUs the go-to choice for an RNN architecture (Goodfellow et al., 2016). The theoretical reason of their success is that their gating mechanism helps to reduce the Vanishing/Exploding gradient problem, and thus to learn long-term dependencies (Hochreiter and Schmidhuber, 1997, Appendix 2).

The names of the different quantities computed in Definition 5.1.11 reflect their intuitive interpretations. The input, forget, and output vectors are all gates: they control the information which will be added to the memory cell based on the new input, the information which will be retained or forgotten from the previous memory cell, and the information which will be transferred from the memory cell to the hidden state, respectively. Notice the identical nature in which all three gates are computed: they are non-linear transformations of affine transformations of the previous hidden state and the input symbol representations. Their parameter matrices define the way in which the gates will influence the memorization, forgetting, and addition of information. The additional information added to the memory cell in the form of the candidate vector \mathbf{g}_t is computed similarly, with the only difference being the activation function. This is the step that bears the most resemblance to the update step of the vanilla RNN (Eq. (5.26)). However, compared to the latter, only parts of this transformation are kept (based on the input and forget vectors \mathbf{i}_t and \mathbf{f}_t). The memory cell \mathbf{c}_t then combines the old memory content \mathbf{c}_{t-1} with the newly integrated information in \mathbf{g}_t to form the new memory content, which is then transformed using the tanh function and combined with the output gate to produce the hidden state \mathbf{h}_t . This is pictorially presented in Fig. 5.3.

As mentioned, the LSTM update step is noticeably more computationally complex than that of a vanilla RNN. This has led to a line of work trying to combine the efficiency of vanilla RNNs and the empirical performance of gated RNNs. In the next subsection, we consider Gated Recurrent Units, one of the best-known compromises found in this domain.

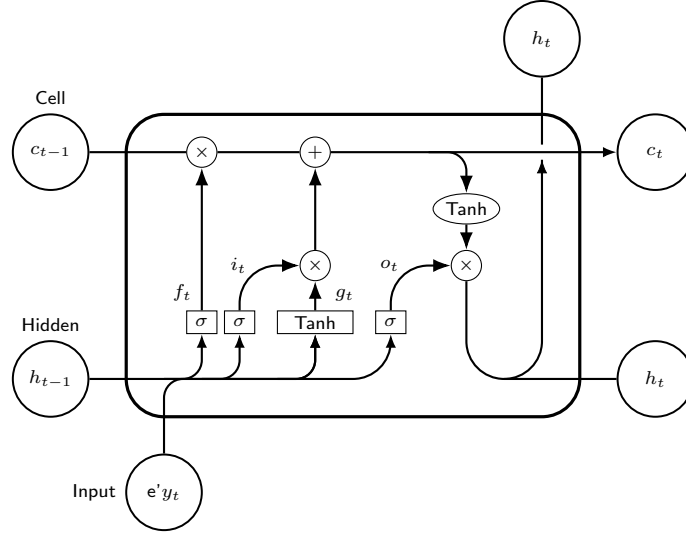


Figure 5.3: A pictorial depiction of the LSTM cell in action. The input i_t , forget f_t , and output o_t gates control which information of the input and of the previous hidden state is retained in the memory cell, and which information is passed to next the hidden state.

Gated Recurrent Units

The Gated Recurrent Unit (GRU, [Choi et al., 2014b,a](#)) provides a compromise between the simplicity of vanilla recurrent neural networks and the empirical success of being able to model long-term dependencies with LSTMs. It defines a gated recurrent update unit that implements a simpler dynamics map by removing the memory component c_t in the LSTM cell and combining the input and forget gates i_t, f_t into one update gate. These changes make GRU more memory efficient and easier to train than LSTM in practice. The full GRU update step is defined as follows.

Definition 5.1.12 (Gated Recurrent Units). *A **gated recurrent unit** defines a dynamics map in which a new hidden state is computed as:*

$$\begin{aligned}
 r_t &= \sigma(\mathbf{U}^r \mathbf{h}_{t-1} + \mathbf{V}^r e'y_t + \mathbf{b}^r) && \text{(reset gate)} \\
 z_t &= \sigma(\mathbf{U}^z \mathbf{h}_{t-1} + \mathbf{V}^z e'y_t + \mathbf{b}^z) && \text{(update gate)} \\
 g_t &= \tanh(\mathbf{U}^g (r_t \odot \mathbf{h}_{t-1}) + \mathbf{V}^g e'y_t + \mathbf{b}^g) && \text{(candidate vector)} \\
 \mathbf{h}_t &= (1 - z_t) \odot g_t + z_t \odot \mathbf{h}_{t-1}
 \end{aligned}$$

r_t, z_t are known as the **reset** and **update** gates, and g_t as the **candidate** vector.

Intuitively, the update gate works like a hot/cold water mixing valve: it is trained to find the *optimum* blend of information of the candidate vector with that coming from the previous hidden state. The reset gate instead, can zero the information of the previous hidden state, when computing the candidate vector. This allows to *forget* past information that becomes irrelevant, exactly like in the LSTM architecture.

Parallelizability: The Achilles' Heel of Recurrent Neural Networks

It is easy to see from the definition of the RNN hidden state (cf. Definition 5.1.3) that, to compute \mathbf{h}_t , we have to compute $\mathbf{h}_{t'}$ for all $t' < t$ first. Another way to say this is that RNNs are inherently *sequential* models, processing the input string one symbol at a time to update their hidden state \mathbf{h}_t , and using this hidden state in turn to compute \mathbf{h}_{t+1} . This results in perhaps the biggest shortcoming of the architecture for its applicability to real-world language modeling: The inability to efficiently *parallelize* the processing (encoding) of long strings. Let us first consider what we mean by the parallelizability of a language model architecture.¹⁰

¹⁰This section provides a very brief and intuitive treatment of parallelization. Our main goal is simply to point out this shortcoming of RNN LMs and with it motivate the next neural architecture we will introduce: transformers.

Due to this sequential nature, the training procedure of RNNs is difficult to *parallelize* effectively, leading to slower training times. This characteristic poses a significant challenge when modeling long strings, as the computation for each element is dependent on the computation of the previous element, leading to a bottleneck in the training process.

In short, in our specific use case of language modeling, parallelization refers to the division of the processing of a specific string across multiple computational nodes, such that any specific node only performs a subset of operations required to process the entire string—the results of the subsets of the operations are then combined to build the representation of the entire string. Importantly, the computations should be performed *independently* between nodes in the sense that no node has to wait for any other node to provide it the results of its computation. Being able to parallelize large models across computational nodes has led to some of the biggest advancements in modern deep learning. As such, parallelizability is a crucial feature of any successful deep learning architecture.

However, notice that any dependence between the computations performed by the nodes defeats the purpose of parallelization—if the nodes have to wait for each other to finish computations, the same operations might as well be performed by a single node. This is where recurrent neural networks fall short: the computations required to encode a string \mathbf{y} into the hidden state $\mathbf{h}_{|\mathbf{y}|}$ will always be sequential, preventing their distribution across different nodes.

Parallelizability in language modeling. When talking about parallelizing language models, it is important to think about *which parts* can actually be parallelized. In the case of RNNs, we saw that no part of the processing can be (besides the matrix multiplication in a single update rule)—the length of the longest chain of dependent computation will always scale linearly with the length of the string. In the next section, we introduce transformers, a recent neural network architecture first introduced for processing text. One of the big contributions of transformers is their parallelizability *during training*—it enables their training on extremely large corpora and is thus one of the main reasons that they are behind the success of many of the most successful modern large language models. Parallelizability *during training* is crucial—notice that parallelization is, in fact, *not possible* during generation from a locally normalized language model (cf. Definition 2.4.5)—by definition, such models will generate *one symbol at a time*. To compute the representation of the new sentence (or the new prefix), which is required for the generation of the next symbol, the generated (sampled) symbol has to first be determined, which leaves their generation process inherently sequential. In that respect, RNNs are as parallelizable as they can be during generation. However, the sequential computation of the hidden states prevents the parallelization of computations of $\text{enc}_{\mathcal{R}}(\mathbf{y}_{\leq t})$ even if the whole string is already *given*.

We will see that the big parallelizability improvements of other architectures only come into play during training when the model is given the whole string in advance (such that no part of the string has to be sequentially generated) and can compute the (log-)likelihood of the given *ground-truth* next symbols given the context. That is, during training and given a string $\mathbf{y} \in \Sigma^*$, the model simply has to compute $p_{\text{SM}}(\bar{y}_t \mid \mathbf{y}_{<t})$ for all $t = 1, \dots, T$ —in representation based language models (cf. Definition 3.1.11) this depends on $\text{enc}(\mathbf{y}_{<t})$. Crucially, computing $p_{\text{SM}}(\bar{y}_t \mid \mathbf{y}_{<t})$ is all that we need for training a language model (in the simplest case that we analyze)—the computed log-likelihood of the ground-truth next character \bar{y}_t is used for computing the loss function during training and used for gradient-based updates to the parameters as discussed in §3.2.3. In the next section, we will see how $\text{enc}(\mathbf{y}_{<t})$ can be computed *without* sequential dependencies. However, in the case of RNNs, $\text{enc}(\mathbf{y}_{<t})$ can only be computed sequentially—even if the entire string is known in advance. This results in a crucial bottleneck in training RNNs on large corpora and vastly limits their applicability to implementing large language models.

5.1.7 Computational Expressiveness of Recurrent Neural Networks

Recurrent neural networks are one of the fundamental and most successful neural language model architectures. In this section, we study some theoretical explanations behind their successes as well as some of their theoretical limitations. Answering this question is essential whenever we require formal guarantees of the correctness of the outputs generated by an LM. For example, one might ask a language model to solve a mathematical problem based on a textual description (Shridhar et al., 2023) or ask it to find an optimal solution to an everyday optimization problem (Lin et al., 2021b). If such problems fall outside the theoretical capabilities

of the LM, we have no ground to believe that the result provided by the model is correct. The question also follows a long line of work on the linguistic capabilities of LMs, as LMs must be able to implement mechanisms of recognizing specific syntactic structures to generate grammatical sequences (Talmor et al., 2020; Hewitt and Manning, 2019; Jawahar et al., 2019; Liu et al., 2019; Icard, 2020a; Manning et al., 2020; Rogers et al., 2021; Belinkov, 2022; Del’etang et al., 2022, *inter alia*).

One way of quantifying the expressive power of computational models is with the complexity of formal languages they can recognize (Del’etang et al., 2022)—we, too, will study the classes of (weighted) formal languages (such as the regular languages and the Turing computable languages) they can express. Through this, diverse formal properties of modern LM architectures have been shown (e.g., Siegelmann and Sontag, 1992; Hao et al., 2018; Korsky and Berwick, 2019; Merrill, 2019; Merrill et al., 2020; Hewitt et al., 2020; Merrill et al., 2022a,b, *inter alia*). Inspecting complex models such as recurrent neural networks through the lens of formal language theory allows us to apply the well-studied theoretical results and understanding from the field to the recently more successful neural models. While studying neural language models, we will revisit various aspects of the classical language models introduced in Chapter 4—indeed, this was also our main motivation for studying those closely.

Specifically, we will focus mainly on the *Elman* recurrent neural networks due to their simplicity and their role as the “fundamental” RNN, capturing their recurrent nature. We will also briefly touch upon the computational power of LSTM networks due to their somewhat different theoretical properties. However, note that most of the results presented in the section generalize to other architectures as well. We begin by investigating Elman RNNs in a practical setting, that is, under relatively strict and realistic assumptions, such as fixed-point arithmetic. We show that Elman RNNs under such a regime are in fact equivalent to weighted finite-state automata. Next, in the second subsection, we show that under more permissive assumptions of infinite precision and unbounded computation time, Elman RNNs are Turing complete.

RNNs and Weighted Regular Languages

Analyzing complex systems with intricate interactions between inputs and parameters and temporal dependencies can be tricky. This is a common issue when studying neural networks in general. In fact, most, if not all, theoretical frameworks for analyzing neural models such as RNNs rely on various assumptions about their components to make the analysis feasible. For example, theoretical results on neural networks (for example, optimization guarantees or function/system identifiability guarantees) often make the assumption that the activation functions are linear or of some other easy-to-analyze form. Similarly, a fruitful manner to analyze the expressivity of recurrent neural networks specifically is by making (somewhat different) simplifying assumptions on the non-linear activation functions, since those are what often make analysis difficult. A common simplification is the use of the Heaviside activation function.

Definition 5.1.13. *The **Heaviside** function is defined as*

$$H(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases} \quad (5.42)$$

In words, the Heaviside function maps every real value either 0 or 1, depending on whether it is greater than or less than zero. In the following, we will refer to the set $\{0, 1\}$ as $\mathbb{B} \stackrel{\text{def}}{=} \{0, 1\}$.

Definition 5.1.14 (Heaviside Elman Network). *A **Heaviside Elman network** (HRNN) is an Elman network with Heaviside function H as the non-linearity.*

Elman network parameters. Importantly, note that the *parameters* of the network do not have to be elements of \mathbb{B} —we assume those can take arbitrary real (or rational) values. Indeed, networks constrained to parameters $\theta \in \mathbb{B}$ would only be able to recognize unweighted languages. Furthermore, for this section, we expand our definition of a real- or rational-weighted RNN to be able to contain *weights* ∞ and $-\infty$. While those are not real (or rational) numbers, we will see they become useful when we want to explicitly exclude specific sequences from the support of the model, i.e., when we want to assign probability 0 to them.

Before we move to the central result of the subsection, we first introduce a fact that makes it easier to talk about how an RNN language models can simulate a deterministic PFSA. We will be interested in conjoining elements of vectors in \mathbb{B}^D , which can be performed by an Elman RNN with appropriately set parameters.

Fact 5.1.4. *Let $m \in [D]$, $i_1, \dots, i_m \in [D]$, and $\mathbf{x}, \mathbf{v} \in \mathbb{B}^D$ with $v_i = \mathbb{1}\{i \in \{i_1, \dots, i_m\}\}$. Then, $H(\mathbf{v}^\top \mathbf{x} - (m-1)) = 1$ if and only if $x_{i_k} = 1$ for all $k = 1, \dots, m$. In other words, $H(\mathbf{v}^\top \mathbf{x} - (m-1)) = x_{i_1} \wedge \dots \wedge x_{i_m}$.*

The central result. The central result of this section is captured in the following theorem.

Theorem 5.1.4 (Equivalence of Heaviside Elman RNNs and WFSAs). *Heaviside Elman RNNs are equivalent to deterministic probabilistic finite-state automata.*

Notice that we only make the claim for *probabilistic* WFSAs. This is without loss of generality, as, from Theorem 4.1.1, we know we can assume \mathcal{A} is locally normalized. We will prove Theorem 5.1.4 by showing that an RNN with Heaviside activations is at most regular, and then showing how such an RNN can in fact simulate any deterministic PFSA. We show each direction as its own lemma.

Lemma 5.1.2. *The distribution represented by a recurrent neural network with a Heaviside non-linearity H is regular.*

Proof. Let $\mathcal{R} = (\Sigma, D, \mathbf{U}, \mathbf{V}, \mathbf{E}, \mathbf{b}_h, \mathbf{h}_0)$ be a HRNN defining the conditional probabilities p_{SM} . We construct a deterministic PFSA $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ defining the same string probabilities. Let $s : \mathbb{B}^D \rightarrow \mathbb{Z}_{2^D}$ be a bijection. Now, for every state $q \stackrel{\text{def}}{=} s(\mathbf{h}) \in Q \stackrel{\text{def}}{=} \mathbb{B}^D$, construct a transition $q \xrightarrow{y/w} q'$ where $q' = \sigma(\mathbf{U}\mathbf{h} + \mathbf{V}\llbracket y \rrbracket + \mathbf{b}_h)$ with the weight $w = p_{\text{SM}}(y \mid \mathbf{h}) = \mathbf{f}_{\Delta[\Sigma]^{-1}}(\mathbf{E}\mathbf{h})_y$. We define the initial function as $\lambda(s(\mathbf{h})) = \mathbb{1}\{\mathbf{h} = \mathbf{h}_0\}$ and final function ρ with $\rho(q) \stackrel{\text{def}}{=} p_{\text{SM}}(\text{EOS} \mid s(q))$.

It is easy to see that \mathcal{A} defined this way is deterministic. We now prove that the weights assigned to strings by \mathcal{A} and \mathcal{R} are the same. Let $\mathbf{y} \in \Sigma^*$ with $|\mathbf{y}| = T$ and

$$\pi = \left(s(\mathbf{h}_0) \xrightarrow{y_1/w_1} q_1, \dots, q_T \xrightarrow{y_T/w_T} q_{T+1} \right)$$

the \mathbf{y} -labeled path starting in $s(\mathbf{h}_0)$ (such a path exists since we the defined automaton is *complete*—all possible transitions are defined for all states).

$$\begin{aligned} \mathcal{A}(\mathbf{y}) &= \lambda(s(\mathbf{h}_0)) \cdot \left[\prod_{t=1}^T w_t \right] \cdot \rho(q_{T+1}) \\ &= 1 \cdot \prod_{t=1}^T p_{\text{SM}}(y_t \mid s^{-1}(q_t)) \cdot p_{\text{SM}}(\text{EOS} \mid s^{-1}(q_{T+1})) \\ &= p_{\text{LN}}(\mathbf{y}) \end{aligned}$$

which is exactly the weight assigned to \mathbf{y} by \mathcal{R} . Note that all paths not starting in $s(\mathbf{h}_0)$ have weight 0 due to the definition of the initial function. \blacksquare

Let us look at an example of the construction above.

Example 5.1.12. *Let $\mathcal{R} = (\Sigma, D, \mathbf{f}, \mathbf{E}, \mathbf{h}_0)$ be a Heaviside RNN sequence model with the parameters*

$$\Sigma = \{a, b\} \tag{5.43}$$

$$D = 2 \tag{5.44}$$

$$\mathbf{f}(\mathbf{h}_t, y) = H \left(\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \mathbf{h}_{t-1} + \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \llbracket y \rrbracket \right) \tag{5.45}$$

$$\mathbf{E} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & -\infty \end{pmatrix} \tag{5.46}$$

$$\mathbf{h}_0 = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \tag{5.47}$$

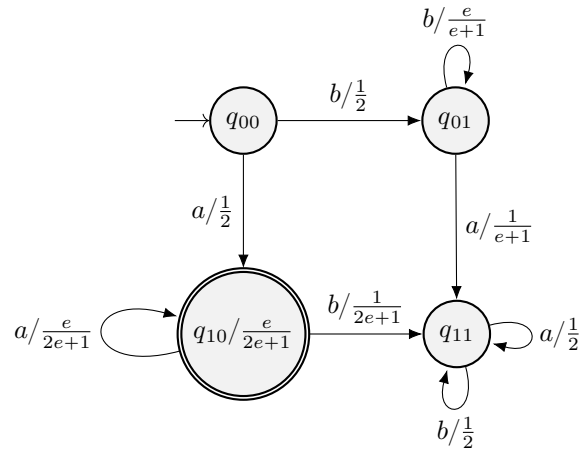


Figure 5.4: The WFSA corresponding to the RNN defined in Eq. (5.43).

and $n(a) = 1$, $n(b) = 2$, and $n(\text{EOS}) = 3$. The automaton corresponding to this RNN contains the states q_{ij} corresponding to the hidden states $\mathbf{h} = \begin{pmatrix} i \\ j \end{pmatrix}$. It is shown in Fig. 5.4; as we can see, the automaton indeed has an exponential number of useful states in the dimensionality of the hidden state, meaning that the RNN is a very compact way of representing it.

To show the other direction of Theorem 5.1.4, we now give a variant of a classic theorem originally due to Minsky (1986) but with a *probabilistic* twist, allowing us to model *weighted* languages with Elman RNNs.

Lemma 5.1.3 (Elman RNNs can encode PFSA). *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a tight probabilistic deterministic finite-state automaton. Then, there exists a Heaviside-activated Elman network with a hidden state of size $\mathbf{h} = |\Sigma||Q|$ that encodes the same distribution as \mathcal{A} .*

We give proof by construction: Given a deterministic PFSA $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$, we construct an Elman RNN $\mathcal{R} = (\Sigma, D, \mathbf{U}, \mathbf{V}, \mathbf{E}, \mathbf{b}_h, \mathbf{h}_0)$ accepting the same weighted language as \mathcal{A} : $L(\mathcal{A}) = L(\mathcal{R})$ by defining the elements of the tuple $(\Sigma, D, \mathbf{U}, \mathbf{V}, \mathbf{E}, \mathbf{b}_h, \mathbf{h}_0)$. In the rest of the section, we will first intuitively describe the construction, and then formally prove the central results that the construction relies on. Let $n: Q \times \Sigma \rightarrow \mathbb{Z}_{|Q||\Sigma|}$ be a bijection, i.e., an ordering of $Q \times \Sigma$, $m: \Sigma \rightarrow \mathbb{Z}_{|\Sigma|}$ an ordering of Σ , and $\bar{m}: \bar{\Sigma} \rightarrow \mathbb{Z}_{|\bar{\Sigma}|}$ an ordering of $\bar{\Sigma}$; these mappings assign each element in their pre-image an integer which can then be used to index into matrices and vectors as we will see below. We use n , m , and \bar{m} to define the one-hot encodings $\llbracket \cdot \rrbracket$ of state-symbol pairs and of the symbols. That is, we assume that $\llbracket q, y \rrbracket_d = \mathbb{1}\{d = n(q, y)\}$, and similar for $\llbracket y \rrbracket$. Similarly to the proof of Lemma 5.1.2, we denote with h and h^{-1} the mappings from $Q \times \bar{\Sigma}$ to the hidden state space of the RNN and its inverse. The alphabet of the RNN of course matches the one of the WFSA.

HRNN’s hidden states. The hidden states of the RNN live in $\mathbb{B}^{|Q||\Sigma|}$. A hidden state \mathbf{h}_t encodes the state q_t the simulated \mathcal{A} is in at time t and the transition symbol y_t with which \mathcal{A} “arrived” at q_t as a *one-hot encoding* of the pair (q_t, y_t) . Formally,

$$\mathbf{h}_t = \llbracket (q_t, y_t) \rrbracket \in \mathbb{B}^{|Q||\Sigma|}. \quad (5.48)$$

This also means that $D = |Q||\Sigma|$. There is a small caveat: how do we set the incoming symbol of \mathcal{A} ’s (sole) initial state q_i (the first time it is entered)? A straightforward solution would be to augment the alphabet of the RNN with the BOS symbol (cf. §2.4), which we define to be the label of the incoming arc denoting the initial state (this would be the only transition labeled with BOS). However, as we show later, the symbol used to arrive into p does not have an effect on the *subsequent* transitions—it is only needed to determine the *target* of the current transition. Therefore, we can simply represent the initial state \mathbf{h}_0 of \mathcal{R} with the one-hot encoding of *any* pair (q_i, a) , where q_i is the initial state of the WFSa and $a \in \Sigma$.

For example, for the fragment of a WFSa in Fig. 5.5, the hidden state encoding the current state q and the incoming arc b is of the form presented in Eq. (5.49).

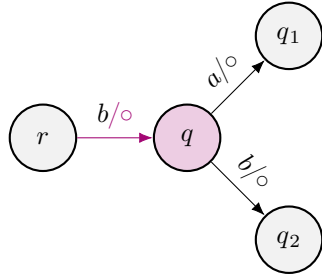


Figure 5.5: A fragment of a WFSA.

$$\mathbf{h}_t = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \leftarrow n(q, b) \quad (5.49)$$

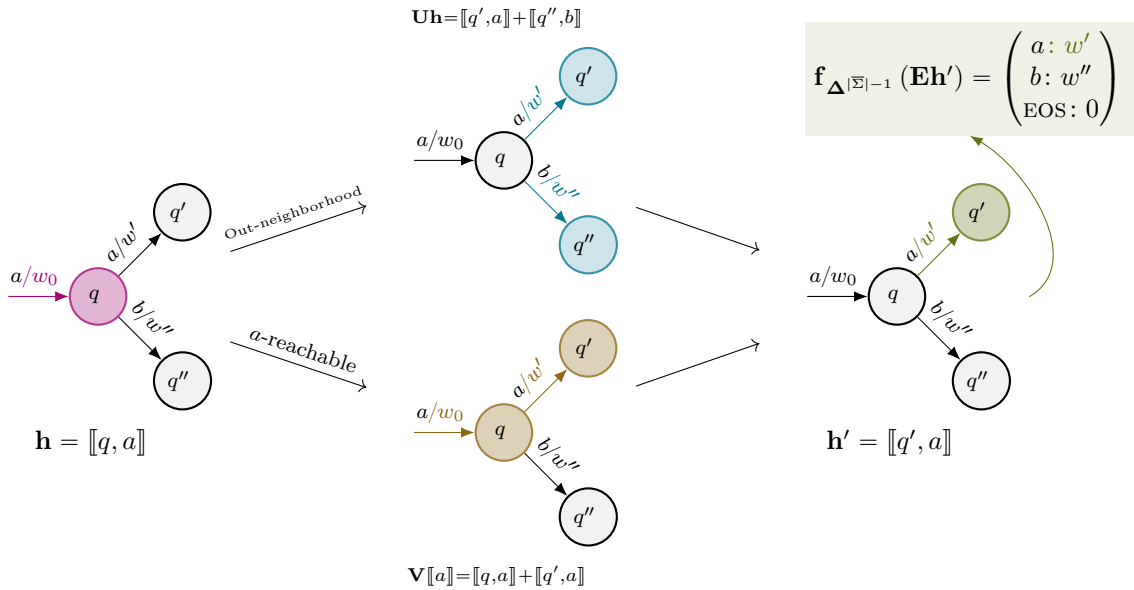


Figure 5.6: A high-level illustration of how the transition function of the FSA is simulated in Minsky's construction on a fragment of an FSA **starting at** q (encoded in \mathbf{h}) and reading the symbol a . The top path disjoins the representations of the states in the **out-neighborhood** of q , whereas the bottom path disjoins the representations of states **reachable by** an a -transition. The Heaviside activation **conjoins** these two representations into \mathbf{h}' (rightmost fragment). Projecting $E\mathbf{h}'$ results in the vector defining the same probability distribution as the outcoming arcs of q (**green box**).

Encoding the transition function. The idea of defining \mathbf{U} , \mathbf{V} , and \mathbf{h}_h is for the Elman update rule to perform, upon reading y_{t+1} , element-wise conjunction between the representations of the out-neighborhood of q_t and the representation of the states \mathcal{A} can transition into after reading in y_{t+1} from *any state*. The former is encoded in the recurrence matrix \mathbf{U} , which has access to the current hidden state that encodes q_t while the latter is encoded in the input matrix \mathbf{V} , which has access to the one-hot representation of y_{t+1} . Conjugating the entries in those two representations will, due to the determinism of \mathcal{A} , result in a single non-zero entry: one representing the state which can be reached from q_t (1st component) using the symbol y_{t+1} (2nd component); see Fig. 5.6.

The recurrence matrix \mathbf{U} lives in $\mathbb{B}^{|\Sigma||Q| \times |\Sigma||Q|}$. The main idea of the construction is for each column $\mathbf{U}[:, n(q, y)]$ of the matrix to represent the “out-neighborhood” of the state q in the sense that the column contains 1’s at the indices corresponding to the state-symbol pairs (q', y') such that \mathcal{A} transitions from q to q' after reading in the symbol y' . That is, for $q, q' \in Q$ and $y, y' \in \Sigma$, we define

$$U_{n(q', y'), n(q, y)} \stackrel{\text{def}}{=} \mathbb{1} \left\{ q_t \xrightarrow{y'/\circ} q' \in \delta \right\}. \quad (5.50)$$

Since y is free, each column is repeated $|\Sigma|$ -times: once for every $y \in \Sigma$ —this is why, after entering the next state, the symbol used to enter it does not matter anymore and, in the case of the initial state, any incoming symbol can be chosen to represent \mathbf{h}_0 .

For example, for the fragment of a WFSFA in Fig. 5.5, the recurrence matrix would take the form

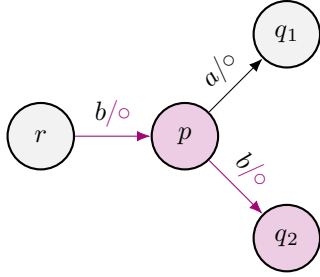
$$\mathbf{U} = \begin{pmatrix} & \overset{n(q, b)}{\downarrow} 0 \\ & \vdots \\ & 1 \\ \dots & \vdots & \dots \\ & 1 \\ & \vdots \\ & 0 \end{pmatrix} \begin{matrix} \leftarrow n(q_1, a) \\ \leftarrow n(q_2, b) \end{matrix} \quad (5.51)$$

and the matrix-vector product $\mathbf{U}\mathbf{h}_t$ with \mathbf{h}_t from before results in

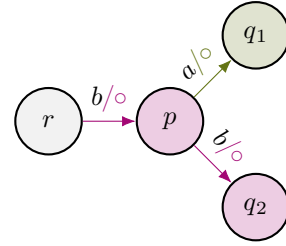
$$\mathbf{U}\mathbf{h}_t = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \begin{matrix} \leftarrow n(q_1, a) \\ \leftarrow n(q_2, b) \end{matrix} \quad (5.52)$$

The input matrix \mathbf{V} lives in $\mathbb{B}^{|\Sigma||Q| \times |\Sigma|}$ and encodes the information about which states can be reached by which symbols (from *any* state in \mathcal{A}). The non-zero entries in the column corresponding to $y' \in \Sigma$ correspond to the state-symbol pairs (q', y') such that q' is reachable with y' from *some* state:

$$V_{n(q', y'), m(y')} \stackrel{\text{def}}{=} \mathbb{1} \left\{ \circ \xrightarrow{y'/\circ} q' \in \delta \right\}. \quad (5.53)$$



(a) An example of a fragment of a WFSA.



(b) An example of a fragment of a WFSA.

For example, for the fragment of a WFSA in Fig. 5.7a, the input matrix would take the form

$$\mathbf{V} = \begin{pmatrix} & \begin{matrix} m(b) \\ \downarrow \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{matrix} & \\ \dots & & \dots \end{pmatrix} \begin{matrix} \leftarrow n(p, b) \\ \leftarrow n(q_2, b) \end{matrix} \quad (5.54)$$

and the matrix-vector product $\mathbf{V}\mathbf{e}_{m(a)}$ and $\mathbf{V}\mathbf{e}_{m(b)}$ would take the form (see also Fig. 5.7b)

$$\mathbf{V}\mathbf{e}_{m(a)} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow n(q_1, a) \quad \mathbf{V}\mathbf{e}_{m(b)} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \begin{matrix} \leftarrow n(p, b) \\ \leftarrow n(q_2, b) \end{matrix} \quad (5.55)$$

Lastly, we define the bias as $\mathbf{b}_h \stackrel{\text{def}}{=} -\mathbf{1} \in \mathbb{R}^{|Q||\Sigma|}$, which allows the Heaviside function to perform the needed conjunction.

To put these components together, consider that, at each step of the computation, \mathcal{R} computes $\mathbf{h}_{t+1} =$

$H(\mathbf{U}\mathbf{h}_t + \mathbf{V}\mathbf{e}_a + \mathbf{b}_h)$ where $y_{t+1} = a$. The input to the non-linearity is computed as follows:

$$\mathbf{U}\mathbf{h}_t + \mathbf{V}\mathbf{e}_{m(a)} + \mathbf{b}_h = \begin{pmatrix} 1 \\ \vdots \\ \color{red}{1} \\ \vdots \\ 1 \\ \vdots \\ 1 \end{pmatrix} \begin{matrix} \leftarrow n(q_1, a) \\ \\ \leftarrow n(q_2, b) \end{matrix} + \begin{pmatrix} 1 \\ \vdots \\ \color{red}{1} \\ \vdots \\ 1 \end{pmatrix} \begin{matrix} \leftarrow n(q_1, a) \end{matrix} + \begin{pmatrix} -1 \\ \vdots \\ -1 \\ \vdots \\ -1 \end{pmatrix} \quad (5.56)$$

The following lemma proves that the construction described correctly implements the transition function of the PFSA.

Lemma 5.1.4. *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a deterministic PFSA, $\mathbf{y} = y_1 \dots y_T \in \Sigma^*$, and q_t the state arrived at by \mathcal{A} upon reading the prefix $\mathbf{y}_{\leq t}$. Let \mathcal{R} be the HRNN specified by the Minsky construction for \mathcal{A} , n the ordering defining the one-hot representations of state-symbol pairs by \mathcal{R} , and \mathbf{h}_t \mathcal{R} 's hidden state after reading $\mathbf{y}_{\leq t}$. Then, it holds that $\mathbf{h}_0 = \llbracket (q_\iota, y) \rrbracket$ where q_ι is the initial state of \mathcal{A} and $y \in \Sigma$ and $\mathbf{h}_T = \llbracket (q_T, y_T) \rrbracket$.*

Proof. Define $s(\mathbf{h}) = \llbracket (q, y) \rrbracket \stackrel{\text{def}}{=} q$. We can then restate the lemma as $s(\mathbf{h}_T) = q_T$ for all $\mathbf{y} \in \Sigma^*$, $|\mathbf{y}| = T$. Let π be the \mathbf{y} -labeled path in \mathcal{A} . We prove the lemma by induction on the string length T .

Base case: $T = 0$. Holds by the construction of \mathbf{h}_0 .

Inductive step: $T > 0$. Let $\mathbf{y} \in \Sigma^*$ with $|\mathbf{y}| = T$ and assume that $s(\mathbf{h}_{T-1}) = q_{T-1}$.

We prove that the specifications of \mathbf{U} , \mathbf{V} , and \mathbf{b}_h ensure that $s(\mathbf{h}_T) = q_T$. By definition of the recurrence matrix \mathbf{U} (cf. Eq. (5.50)), the vector $\mathbf{U}\mathbf{h}_{T-1}$ will contain a 1 at the entries $n(q', y')$ for $q' \in Q$ and $y' \in \Sigma$ such that $q_T \xrightarrow{y'/\circ} q' \in \delta$. This can equivalently be written as $\mathbf{U}\mathbf{h}_{T-1} = \bigvee_{q_T \xrightarrow{y'/\circ} q' \in \delta} \llbracket (q', y') \rrbracket$, where the disjunction is applied element-wise.

On the other hand, by definition of the input matrix \mathbf{V} (cf. Eq. (5.53)), the vector $\mathbf{V}\llbracket y_T \rrbracket$ will contain a 1 at the entries $n(q', y_T)$ for $q' \in Q$ such that $\circ \xrightarrow{y_T/\circ} q' \in \delta$. This can also be written as $\mathbf{V}\llbracket y_T \rrbracket = \bigvee_{\circ \xrightarrow{y_T/\circ} q' \in \delta} \llbracket (q', y_T) \rrbracket$.

By Fact 5.1.4, $H(\mathbf{U}\mathbf{h}_{T-1} + \mathbf{V}\llbracket y_T \rrbracket + \mathbf{b}_h)_{n(q', y')} = H(\mathbf{U}\mathbf{h}_{T-1} + \mathbf{V}\llbracket y_T \rrbracket - \mathbf{1})_{n(q', y')} = 1$ holds if and only if $(\mathbf{U}\mathbf{h}_{T-1})_{n(q', y')} = 1$ and $(\mathbf{V}\llbracket y_T \rrbracket)_{n(q', y')} = 1$. This happens if

$$q_T \xrightarrow{y'/\circ} q' \in \delta \text{ and } \circ \xrightarrow{y_T/\circ} q' \in \delta \iff q_T \xrightarrow{y_T/\circ} q', \quad (5.57)$$

i.e., if and only if \mathcal{A} transitions from q_T to q_T upon reading y_T (it transitions only to q_T due to determinism).

Since the string \mathbf{y} was arbitrary, this finishes the proof. \blacksquare

Encoding the transition probabilities. We now turn to the second part of the construction: encoding the string acceptance weights given by \mathcal{A} into the probability distribution defined by \mathcal{R} . We present two ways of doing that: using the more standard softmax formulation, where we make use of the extended real numbers, and with the sparsemax.

The conditional probabilities assigned by \mathcal{R} are controlled by the $|\bar{\Sigma}| \times |Q||\Sigma|$ -dimensional output matrix \mathbf{E} . Since \mathbf{h}_t is a one-hot encoding of the state-symbol pair q_t, y_t , the matrix-vector product $\mathbf{E}\mathbf{h}_t$ simply looks up the values in the $n(q_t, y_t)^{\text{th}}$ column. After being projected to $\Delta^{|\bar{\Sigma}|-1}$, the entry in the projected vector corresponding to some $y_{t+1} \in \bar{\Sigma}$ should match the probability of that symbol given that \mathcal{A} is in the state q_t .

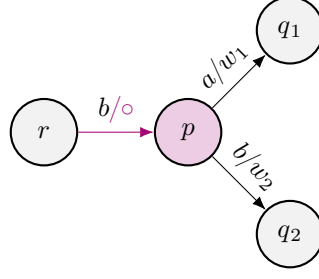


Figure 5.8: An example of a fragment of a WFSA.

This is easy to achieve by simply encoding the weights of the outgoing transitions into the $n(q_t, y_t)^{\text{th}}$ column, depending on the projection function used. This is especially simple in the case of the sparsemax formulation. By definition, in a PFSA, the weights of the outgoing transitions and the final weight of a state q_t form a probability distribution over $\bar{\Sigma}$ for every $q_t \in Q$. Projecting those values to the probability simplex, therefore, leaves them intact. We can therefore define

$$\mathbf{E}_{\bar{m}(y')n(q,y)} \stackrel{\text{def}}{=} \begin{cases} \omega(q \xrightarrow{y'/w} \circ) & | \text{ if } y' \in \Sigma \\ \rho(q) & | \text{ otherwise} \end{cases} \quad (5.58)$$

Projecting the resulting vector $\mathbf{E}\mathbf{h}_t$, therefore, results in a vector whose entries represent the transition probabilities of the symbols in $\bar{\Sigma}$.

In the more standard softmax formulation, we proceed similarly but log the non-zero transition weights. Defining $\log 0 \stackrel{\text{def}}{=} -\infty$,¹¹ we set

$$\mathbf{E}_{\bar{m}(y')n(q,y)} \stackrel{\text{def}}{=} \begin{cases} \log \omega(q \xrightarrow{y'/w} \circ) & | \text{ if } y' \in \Sigma \\ \log \rho(q) & | \text{ otherwise} \end{cases} \quad (5.59)$$

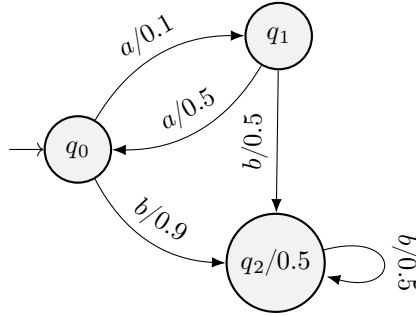
It is easy to see that the entries of the vector $\text{softmax}(\mathbf{E}\mathbf{h}_t)$ form the same probability distribution as the original outgoing transitions out of q . Over the course of an entire input string, these weights are multiplied as the RNN transitions between different hidden states corresponding to the transitions in the original PFSA \mathcal{A} .

For example, for the fragment of a WFSA in Fig. 5.8, the output matrix would take the form

$$\mathbf{E} = \begin{pmatrix} & n(q, b) \\ & \downarrow \\ & -\infty \\ & \vdots \\ & \log w_1 \\ \cdots & \vdots & \cdots \\ & \log w_2 \\ & \vdots \\ & -\infty \end{pmatrix} \begin{matrix} \leftarrow \bar{m}(a) \\ \leftarrow \bar{m}(b) \end{matrix} \quad (5.60)$$

This means that, if \mathbf{h}_t encodes the state-symbol pair (q, y) , the vector $\mathbf{E}\mathbf{h}_t$ will copy the selected column in \mathbf{E} which contains the output weight for all out symbols y ; of q , i.e., the entry $\mathbf{E}\mathbf{h}_{\bar{m}(y')}$ contains the weight on the arc $q \xrightarrow{y'/w} \circ$. Over the course of an entire input string \mathbf{y} , these probabilities are simply multiplied as the RNN transitions between different hidden states corresponding to the transitions in the original WFSA \mathcal{A} .

¹¹Note that the $-\infty$ entries are only needed whenever the original WFSA assigns 0 probability to some transitions. In many implementations using softmax-activated probabilities, this would not be required.

Figure 5.9: The WFS \mathcal{A} .

For example, for the fragment of a WFS \mathcal{A} in Fig. 5.8, the matrix-vector product $\mathbf{E}\mathbf{h}_t$ would take the form

$$\mathbf{E}\mathbf{h}_t = \begin{pmatrix} -\infty \\ \vdots \\ \log w_1 \\ \vdots \\ \log w_2 \\ \vdots \\ -\infty \end{pmatrix} \begin{matrix} \leftarrow \overline{m}(a) \\ \\ \leftarrow \overline{m}(b) \end{matrix} \quad (5.61)$$

The equivalence of the produced RNN LM to the PFSA is shown in the following lemma.

Lemma 5.1.5. *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a deterministic PFSA, $\mathbf{y} = y_1 \dots y_T \in \Sigma^*$, and q_t the state arrived at by \mathcal{A} upon reading the prefix $\mathbf{y}_{\leq t}$. Let \mathcal{R} be the HRNN specified by the Minsky construction for \mathcal{A} , \mathbf{E} the output matrix specified by the generalized Minsky construction, n the ordering defining the one-hot representations of state-symbol pairs by \mathcal{R} , and \mathbf{h}_t \mathcal{R} 's hidden state after reading $\mathbf{y}_{\leq t}$. Then, it holds that $p_{\text{LN}}(\mathbf{y}) = \mathcal{A}(\mathbf{y})$.*

Proof. Let $\mathbf{y} \in \Sigma^*$, $|\mathbf{y}| = T$ and let π be the \mathbf{y} -labeled path in \mathcal{A} . Again, let $\bar{p}(\mathbf{y}) \stackrel{\text{def}}{=} \prod_{t=1}^{|\mathbf{y}|} p_{\text{SM}}(y_t \mid \mathbf{y}_{<t})$. We prove $\bar{p}(\mathbf{y}) = \prod_{t=1}^T w_t$ by induction on T .

Base case: $T = 0$. In this case, $\mathbf{y} = \varepsilon$, i.e., the empty string, and $\mathcal{A}(\varepsilon) = 1$. \mathcal{R} computes $\bar{p}(\varepsilon) = \prod_{t=1}^0 p_{\text{SM}}(y_t \mid \mathbf{y}_{<t}) = 1$.

Inductive step: $T > 0$. Assume that the $\bar{p}(y_1 \dots y_{T-1}) = \prod_{t=1}^{T-1} w_t$. By Lemma 5.1.4, we know that $s(\mathbf{h}_{T-1}) = q_T$ and $s(\mathbf{h}_T) = q_T$. By the definition of \mathbf{E} for the specific $\mathbf{f}_{\Delta|\Sigma|-1}$, it holds that $\mathbf{f}_{\Delta|\Sigma|-1}(\mathbf{E}\mathbf{h}_{T-1})_{m(\mathbf{y})} = \omega(s(\mathbf{h}_{T-1}) \xrightarrow{y/w_T} s(\mathbf{h}_T)) = w_T$. This means that $\bar{p}(\mathbf{y}_{\leq T}) = \prod_{t=1}^T w_t$, which is what we wanted to prove.

Clearly, $p_{\text{LN}}(\mathbf{y}) = \bar{p}(\mathbf{y}) p_{\text{SM}}(\text{EOS} \mid \mathbf{y})$. By the definition of \mathbf{E} (cf. Eq. (5.58)), $(\mathbf{E}\mathbf{h}_T)_{m(\text{EOS})} = \rho(s(\mathbf{h}_T))$, meaning that

$$p_{\text{LN}}(\mathbf{y}) = \bar{p}(\mathbf{y}) p_{\text{SM}}(\text{EOS} \mid \mathbf{y}) = \prod_{t=1}^T w_t \rho(s(\mathbf{h}_T)) = \mathcal{A}(\mathbf{y}).$$

Since $\mathbf{y} \in \Sigma^*$ was arbitrary, this finishes the proof. ■

We now walk through an example of the Minsky construction.

Example 5.1.13. *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a WFS \mathcal{A} as shown in Fig. 5.9. Since \mathcal{A} has $|Q| = 3$ states and an alphabet of $|\Sigma| = 2$ symbols, the hidden state of the representing RNN \mathcal{R} will be of dimensionality*

$3 \cdot 2 = 6$. Assume that the set of state-symbol pairs is ordered as $(q_0, a), (q_0, b), (q_1, a), (q_1, b), (q_2, a), (q_2, b)$. The initial state can be represented (choosing a as the arbitrary “incoming symbol”) as

$$\mathbf{h}_0 = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}. \quad (5.62)$$

The recurrent matrix \mathbf{U} of \mathcal{R} is

$$\mathbf{U} = \begin{pmatrix} 0 & 0 & \color{red}{1} & \color{red}{1} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \color{red}{1} & \color{red}{1} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \color{red}{1} & \color{red}{1} & \color{red}{1} & \color{red}{1} & \color{red}{1} & \color{red}{1} \end{pmatrix}, \quad (5.63)$$

the input matrix \mathbf{V}

$$\mathbf{V} = \begin{pmatrix} \color{red}{1} & 0 \\ 0 & 0 \\ \color{red}{1} & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & \color{red}{1} \end{pmatrix}, \quad (5.64)$$

and the output matrix \mathbf{E} is

$$\mathbf{E} = \begin{pmatrix} \log(0.1) & \log(0.1) & \log(0.5) & \log(0.5) & -\infty & -\infty \\ \log(0.9) & \log(0.9) & \log(0.5) & \log(0.5) & \log(0.5) & \log(0.5) \\ -\infty & -\infty & -\infty & -\infty & \log(0.5) & \log(0.5) \end{pmatrix}, \quad (5.65)$$

where the last row corresponds to the symbol EOS. The target of the b -labeled transition from q_0 ($q_0 \xrightarrow{b/0.9} q_2$) is computed as follows:

$$\begin{aligned} \mathbf{h}_1 &= H(\mathbf{U}\mathbf{h}_0 + \mathbf{V}[[b]] + \mathbf{b}_h) \\ &= H \left(\begin{pmatrix} \color{green}{0} & 0 & 1 & 1 & 0 & 0 \\ \color{green}{0} & 0 & 0 & 0 & 0 & 0 \\ \color{green}{1} & 1 & 0 & 0 & 0 & 0 \\ \color{green}{0} & 0 & 0 & 0 & 0 & 0 \\ \color{green}{0} & 0 & 0 & 0 & 0 & 0 \\ \color{green}{1} & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} \color{green}{1} \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 1 & \color{green}{0} \\ 0 & \color{green}{0} \\ 1 & \color{green}{0} \\ 0 & \color{green}{0} \\ 0 & \color{green}{0} \\ 0 & \color{green}{1} \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} + \begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \\ -1 \\ -1 \end{pmatrix} \right) \\ &= H \left(\begin{pmatrix} 0 \\ 0 \\ \color{green}{1} \\ 0 \\ 0 \\ \color{green}{1} \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \color{green}{1} \end{pmatrix} + \begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \\ -1 \\ -1 \end{pmatrix} \right) = H \left(\begin{pmatrix} -1 \\ -1 \\ \color{blue}{0} \\ -1 \\ -1 \\ \color{green}{1} \end{pmatrix} \right) = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \color{green}{1} \end{pmatrix}, \end{aligned}$$

which corresponds exactly the configuration in which \mathcal{A} is in state q_2 which it arrived to by reading in the symbol b .

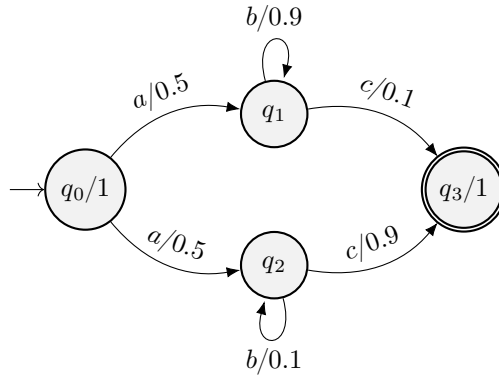


Figure 5.10: A non-determinizable PFSA. It assigns the string $ab^n c$ the probability $\mathcal{A}(ab^n c) = 0.5 \cdot 0.9^n \cdot 0.1 + 0.5 \cdot 0.1^n \cdot 0.9$, which can not be expressed as a single term for arbitrary $n \in \mathbb{N}_{\geq 0}$.

The probability of the string $\mathbf{y} = b$ under the locally-normalized model induced by \mathcal{R} can be computed as

$$\begin{aligned}
 p_{LN}(\mathbf{y}) &= p_{LN}(b) = p_{SM}(b \mid \text{BOS}) p_{SM}(\text{EOS} \mid b) = p_{SM}(b \mid \mathbf{h}_0) p_{SM}(\text{EOS} \mid \mathbf{h}_1) \\
 &= \text{softmax}(\mathbf{E}\mathbf{h}_0)_b \text{softmax}(\mathbf{E}\mathbf{h}_1)_{\text{EOS}} \\
 &= \text{softmax} \left(\begin{pmatrix} \log(0.1) & \cdots & -\infty \\ \log(0.9) & \cdots & \log(0.5) \\ -\infty & \cdots & \log(0.5) \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}_b \right) \\
 &\quad \text{softmax} \left(\begin{pmatrix} \log(0.1) & \cdots & -\infty \\ \log(0.9) & \cdots & \log(0.5) \\ -\infty & \cdots & \log(0.5) \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}_{\text{EOS}} \right) \\
 &= \text{softmax} \begin{pmatrix} \log(0.1) \\ \log(0.9) \\ -\infty \end{pmatrix}_b \cdot \text{softmax} \begin{pmatrix} -\infty \\ \log(0.5) \\ \log(0.5) \end{pmatrix}_{\text{EOS}} = 0.9 \cdot 0.5 = 0.45.
 \end{aligned}$$

Implications for recurrent neural language models. Lemmas 5.1.2 and 5.1.3 formalize the equivalence between HRNNs and deterministic PFSA. A direct corollary of this result is that HRNNs are at most as expressive as deterministic PFSA and, therefore, strictly less expressive as general, non-deterministic, PFSA.¹² An example of a very simple non-deterministic PFSA, i.e., a PFSA whose distribution cannot be expressed by an HRNN LM, is shown in Fig. 5.10. Furthermore, even if a non-deterministic PFSA can be determinized, the number of states of the determinized machine can be exponential in the size of the non-deterministic one (Buchsbaum et al., 2000). In this sense, non-deterministic PFSA can be seen as exponentially compressed representations of finite-state LMs. However, the compactness of this non-deterministic representation must be “undone” using determinization before it can be encoded by an HRNN.

While Lemma 5.1.2 focuses on HRNN LMs and shows that they are finite-state, a similar argument could be made for any RNN whose activation functions map onto a finite set. This is the case with any implementation of an RNN on a computer with finite-precision arithmetic—in that sense, all deployed RNNLMs are finite-state, albeit very large in the sense of encoding possibly very large weighted finite-state automata. However, there are a few important caveats with this: firstly, notice that, although finite, the

¹²General PFSA are, in turn, equivalent to probabilistic regular grammars and discrete Hidden Markov Models (Icard, 2020b).

number of states represented by an RNN is *exponential* in the size of the hidden state. Even for moderate hidden state dimensionalities, this can be very large (hidden states can easily be of size 100–1000). In other words, one can view RNNs as very compact representations of large deterministic probabilistic finite-state automata whose transition functions are represented by the RNN’s update function. Furthermore, since the topology of this implicit WFSa is completely determined by the update function of the RNN, it can be *learned* very flexibly yet efficiently based on the training data—this is made possible by the *sharing* of parameters across the entire graph of the WFSa instead of explicitly parametrizing every possible transition, as, for example, in §4.1.3, or hard-coding the allowed transitions as in §4.1.5. This means that the WFSa is not only represented, but also *parametrized* very efficiently by an RNN. Nevertheless, there is an important detail that we have somewhat neglected so far: this is the requirement that the simulated WFSa be *deterministic*.

Constructing a smaller RNN: Lower Bounds on the Space Complexity of Simulating PFSA with RNNs

Lemma 5.1.3 shows that HRNN LMs are at least as powerful as dPFSA. More precisely, it shows that any dPFSA $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ can be simulated by an HRNN LM of size $\mathcal{O}(|Q||\Sigma|)$. In this section, we address the following question: How large does an HRNN LM have to be such that it can correctly simulate a dPFSA? We study the asymptotic bounds with respect to the size of the set of states, $|Q|$, as well as the number of symbols, $|\Sigma|$.

Asymptotic Bounds in $|Q|$ Intuitively, the 2^D configurations of a D -dimensional HRNN hidden state could represent 2^D states of a (P)FSA. One could therefore expect that we could achieve exponential compression of a dPFSA by representing it as an HRNN LM. Interestingly, this is not possible in general: extending work by Dewdney (1977), Indyk (1995) shows that, to represent an unweighted FSA with an HRNN, one requires an HRNN of size $\Omega(|\Sigma|\sqrt{|Q|})$. This lower bound can be achieved. For completeness, we present constructions by Dewdney (1977); Indyk (1995), which represent an unweighted FSA with a HRNN of size $\mathcal{O}(|\Sigma||Q|^{\frac{3}{4}})$ and $\mathcal{O}(|\Sigma|\sqrt{|Q|})$, respectively, next, before giving a lower bound in for the probabilistic case.

Lemma 5.1.3 gives a relatively simple construction of an RNN recognizing a weighted regular language. However, the resulting RNN is relatively large, with a hidden state of size linear in the number of states of the (deterministic) WFSa recognizing the language, with the additional multiplicative factor in the size of the alphabet. Note that constructions resulting in smaller RNNs exist, at least for the unweighted case. For example, for an arbitrary WFSa $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$, Dewdney (1977); Alon et al. (1991) present a construction of an RNN with a hidden state of size $\mathcal{O}(|Q|^{\frac{3}{4}})$ simulating \mathcal{A} , whereas Indyk (1995) provides a construction of an RNN with a hidden state of size $\mathcal{O}(|Q|^{\frac{1}{2}})$. The latter is also provably a *lower bound* on the number of neurons required to represent an arbitrary unweighted FSA with a Heaviside-activated recurrent neural network (Indyk, 1995). It is not yet clear if this can be generalized to the weighted case or if Minsky’s construction is indeed optimal in this setting. This is quite interesting since one would expect that an RNN with a hidden state of size D can represent up to 2^D individual states (configurations of the D -dimensional vector). However, the form of the transition function with the linear transformation followed by a Heaviside activation limits the number of transition functions that can be represented using D dimensions, resulting in the required exponential increase in the size of the hidden state.

Minsky’s construction (Lemma 5.1.3) describes how to represent a dPFSA \mathcal{A} with a HRNN of size linear in the number of \mathcal{A} ’s states. Importantly, the encoding of the FSA transition function (taken from Minsky’s original construction) is decoupled from the parameter defining the probability distribution, \mathbf{E} . This section describes two asymptotically more space-efficient ways of constructing the component simulating the transition function. They originate in the work by Dewdney (1977), who showed that an unweighted FSA $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ can be represented by an HRNN of size $\mathcal{O}(|\Sigma||Q|^{\frac{3}{4}})$. Using the same ideas, but a specific trick to compress the size of the processing layer of the RNN further, Indyk (1995) reduced this bound to $\mathcal{O}(|\Sigma|\sqrt{|Q|})$, which, as discussed in §5.1.7, is asymptotically optimal. Naturally, as shown in §5.1.7, the space-efficiency gain can not be carried over to the weighted case—that is, the space-efficiency is asymptotically overtaken by the output matrix \mathbf{E} . Nevertheless, for a more complete treatment of the subject,

we cover the two compressed constructions of the HRNN simulating an unweighted FSA in this section in our notation. Importantly, given a dPFSA, we focus only on the *underlying FSA*, i.e., the unweighted transition function of the automaton, since by Theorem 5.1.7, the compression can only be achieved with components representing that part of the automaton.

Dewdney's Construction

This section describes the construction due to Dewdney (1977) in our notation. Since some of the parts are very similar to the construction due to Indyk (1995), those parts are reused in §5.1.7 and introduced more generally.

Representing states of the FSA. Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be a deterministic FSA. Recall that Minsky's construction encodes the \mathcal{A} 's current state as a one-hot encoding of the state-symbol pair. The construction due to Dewdney (1977), on the other hand, represents the states separately from the symbols. It encodes the states with *two-hot representations* by using the coefficients of what we call a square-root state representation. This results in representations of states of size $\mathcal{O}(\sqrt{|Q|})$. The input symbols are incorporated into the hidden state *separately*.¹³

Definition 5.1.15. Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be an FSA and $s \stackrel{\text{def}}{=} \lceil \sqrt{|Q|} \rceil$. We define the **square-root state representation** of \mathcal{A} 's states $q \in Q$ as¹⁴

$$\phi_2(q) \stackrel{\text{def}}{=} \left(\left\lfloor \frac{q}{s} \right\rfloor, q \bmod s \right). \quad (5.66)$$

We denote the inverse of ϕ_2 with ϕ_2^{-1} and further define for $k \in \mathbb{Z}_s$

$$\phi_2^{-1}(k, \cdot) \stackrel{\text{def}}{=} \{q \in Q \mid \varphi_0 = k \text{ where } \varphi = \phi_2(q)\} \quad (5.67)$$

and $\phi_2^{-1}(\cdot, k)$ analogously.

Specifically, we will denote $\phi_2^{-1}(k, \cdot)$ and $\phi_2^{-1}(\cdot, k)$ with k in the j^{th} position (with $j \in \mathbb{Z}_2$, 0 for $\phi_2^{-1}(k, \cdot)$ and 1 for $\phi_2^{-1}(\cdot, k)$) as $\Phi_{k,j}$.

We can think of the function ϕ_2 as representing states of the FSA in a two-dimensional space $\mathbb{Z}_s \times \mathbb{Z}_s$. However, to efficiently simulate \mathcal{A} with an HRNN, it is helpful to think of $\phi_2(q)$ in two different ways: as a vector $\mathbf{v} \in \mathbb{N}_{\geq 0}^{2|Q|}$, or as a matrix in $\mathbb{B}^{|Q| \times |Q|}$ in the following sense.

Definition 5.1.16. Given a square-root state representation function ϕ_2 , we define the **vector representation** of the state $q \in Q$ as the vector $\mathbf{v}(q) \in \mathbb{B}^{2|Q|}$ with

$$\mathbf{v}(q)_{\varphi_0} = 1 \quad (5.68)$$

$$\mathbf{v}(q)_{s+\varphi_1} = 1, \quad (5.69)$$

where $\varphi = (\varphi_0, \varphi_1) = \phi_2(q)$, and all other entries 0. Furthermore, we define the **matrix representation** of the state $q \in Q$ as the matrix $\mathbf{B} \in \mathbb{B}^{|Q| \times |Q|}$ with

$$\mathbf{W}q_{\varphi_0\varphi_1} = 1 \quad (5.70)$$

and all other entries 0.

Dewdney's construction also heavily relies on the representations of sets of states. We define those additively.

¹³This again adds a factor $|\Sigma|$ to the size of the hidden state, as we discuss later.

¹⁴Notice that $\phi_2(q)$ represents the coefficients of the expression of $q \in \mathbb{N}$ in base s .

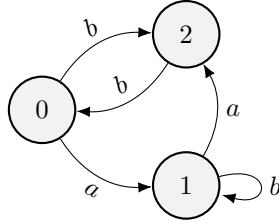


Figure 5.11: An example of a fragment of an FSA.

Definition 5.1.17. Let $\mathcal{Q} \subseteq Q$ be a set of states. We define the vector representation of \mathcal{Q} as the vector

$$\mathbf{v}(\mathcal{Q}) \stackrel{\text{def}}{=} \bigvee_{q \in \mathcal{Q}} \mathbf{v}(q). \quad (5.71)$$

Similarly, we define the matrix representation of \mathcal{Q} as the matrix

$$\mathbf{W}_Q \mathcal{Q} \stackrel{\text{def}}{=} \bigvee_{q \in \mathcal{Q}} \mathbf{W}_Q q. \quad (5.72)$$

To help understand the above definitions, we give an example of a FSA and the representations of its states.

Example 5.1.14. Consider the FSA in Fig. 5.11, for which $s = \lceil \sqrt{|Q|} \rceil = \lceil \sqrt{3} \rceil = 2$, meaning that

$$\phi_2(0) = (0, 0) \quad \phi_2(1) = (0, 1) \quad \phi_2(2) = (1, 0), \quad (5.73)$$

resulting in the state-to-vector mapping¹⁵

$$\mathbf{v}(0) = \begin{pmatrix} 1 & 0 & | & 1 & 0 \end{pmatrix} \quad (5.74)$$

$$\mathbf{v}(1) = \begin{pmatrix} 1 & 0 & | & 0 & 1 \end{pmatrix} \quad (5.75)$$

$$\mathbf{v}(2) = \begin{pmatrix} 0 & 1 & | & 1 & 0 \end{pmatrix}, \quad (5.76)$$

and the state-to-matrix mapping

$$\mathbf{W}_Q 0 = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \quad \mathbf{W}_Q 1 = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \quad \mathbf{W}_Q 2 = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}. \quad (5.77)$$

The two components of the vector representations separated by “|” denote the two halves of the representation vectors, corresponding to the two components of $\phi_2(q)$.

High-level idea of Dewdney’s construction. Given these definitions, the intuition behind Dewdney’s construction of an HRNN simulating an FSA \mathcal{A} is the following:

1. Represent \mathcal{A} ’s states as vectors in \mathbb{B}^{2s} , or, equivalently, matrices in $\mathbb{B}^{s \times s}$.
2. For each $q \in Q$, construct the matrix representation of the set of y -predecessors $\mathbf{W}_Q \text{Pred}(q; y)$ for all $y \in \Sigma$.
3. To simulate \mathcal{A} ’s transition function δ , compare the representation of the current state q_t with all constructed predecessor matrices $\mathbf{W}_Q \text{Pred}(q; y_t)$ given the current input symbol y_t . Activate the two-hot representation of the (unique) state q_{t+1} for which the representation of q_t was detected in q_{t+1} ’s predecessor matrix for symbol y_t , $\mathbf{W}_Q \text{Pred}(q_{t+1}; y_t)$.

¹⁵Despite the notation $(\dots | \dots)$, we assume we are working with column vectors.

Simulating the transition function of an FSA by detecting preceding states. We elaborate on the last point above since it is the central part of the construction.¹⁶ The idea of simulating the transition function δ is reduced to detecting *whose predecessor* given the current input symbol y_t is currently active—naturally, this should be the state active at $t + 1$. Concretely, consider again the FSA \mathcal{A} in Fig. 5.11. The predecessors of the three states, indexed by the incoming symbols are: for 0 $\{b : 2\}$, for 1 $\{a : 1, b : 0\}$, and for 2 $\{a : 1, b : 0\}$. Suppose that at some time t , \mathcal{A} is in state 0 and is reading in the symbol b . Then, since the state 0 is the b -predecessor of the state 2, we know that at time $t + 1$, \mathcal{A} will be in state 2. This principle can be applied more generally: to determine the state of an FSA at time $t + 1$, we simply have to somehow detect whose *predecessor* is active at time t *given the current input symbol* at time t .

The crux of Dewdney’s construction is then the following:¹⁷ How do we, using only the Elman update rule, determine whose y_t -predecessor is active at time t ? This can be done by *detecting* which predecessor matrix $\mathbf{W}_Q \text{Pred}(q; y_t)$ the representation of the current state q_t is included in in the sense that if $\phi_2(q_t) = \varphi$, it holds that $\mathbf{W}_Q \text{Pred}(q; y_t)_{\varphi_0 \varphi_1} = 1$. To be able to formally talk about the detection of a representation in a set of predecessors, we define several notions of **matrix detection**.

Informally, we say that a matrix is easily detectable if the presence of its non-zero elements can be detected using a single neuron in the hidden layer of a HRNN.

Definition 5.1.18. Let $\mathbf{B} \in \mathbb{B}^{D \times D}$ be a binary matrix. We say that \mathbf{B} is **easily detectable** if there exist $\mathbf{w} \in \mathbb{Q}^{2D}$ and $b \in \mathbb{Q}$ (neuron coefficients) such that

$$\sigma(\langle \mathbf{e}_{ij}, \mathbf{w} \rangle + b) = 1 \iff B_{ij} = 1, \quad (5.78)$$

where $\mathbf{e}_{ij} = (\mathbf{e}_i \mid \mathbf{e}_j)$ refers to the $2D$ -dimensional vector with 1’s at positions i and $D + j$. In words, this means that the neuron defined by \mathbf{w}, b fires on the input \mathbf{e}_{ij} if and only if $B_{ij} = 1$.

We define detectable matrices as the matrices which can be detected using a conjunction of two neurons.

Definition 5.1.19. Let $\mathbf{B} \in \mathbb{B}^{D \times D}$ be a binary matrix. We say that \mathbf{B} is **detectable** if there exist $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{Q}^{2D}$ and $b_1, b_2 \in \mathbb{Q}$ such that

$$\sigma(\langle \mathbf{e}_{ij}, \mathbf{w}_1 \rangle + b_1) = 1 \wedge \sigma(\langle \mathbf{e}_{ij}, \mathbf{w}_2 \rangle + b_2) = 1 \iff B_{ij} = 1. \quad (5.79)$$

Furthermore, we say that a matrix is (easily) **permutation-detectable** if there exist permutation matrices \mathbf{P} and \mathbf{Q} such that \mathbf{PBQ} is (easily) detectable.

Intuitively, this means that one can effectively replace an easily detectable matrix \mathbf{B} with a *single neuron*: instead of specifying the matrix explicitly, one can simply detect if an entry B_{ij} of \mathbf{B} is 1 by passing \mathbf{e}_{ij} through the neuron and seeing if it fires. This reduces the space complexity from D^2 to $2D$. Similarly, one can replace a detectable matrix with *two* neurons. As shown in Fact 5.1.4, the required conjunction of the two resulting neurons can then easily be performed by a third (small) neuron, meaning that a detectable matrix is effectively represented by a two-layer MLP.

An example of easily detectable matrices are the so-called **northwestern** matrices.

Definition 5.1.20. A matrix $\mathbf{B} \in \mathbb{B}^{D \times D}$ is **northwestern** if there exists a vector α with $|\alpha| = D$ and $D \geq \alpha_1 \geq \dots \geq \alpha_D \geq 0$ such that

$$B_{ij} = 1 \iff j \leq \alpha_i. \quad (5.80)$$

Intuitively, northwestern matrices contain all their ones contiguously in their upper left (northwest) corner. An example of a northwestern matrix for $\alpha = (2 \ 1 \ 1)$ is

$$\mathbf{B} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}. \quad (5.81)$$

Lemma 5.1.6. Northwestern matrices are easily detectable.

¹⁶Later, we will see that Indyk (1995) uses the exact same idea for simulating δ .

¹⁷Again, the same applies to Indyk (1995).

Proof. Define

$$\mathbf{w} \stackrel{\text{def}}{=} (\boldsymbol{\alpha} \mid D \ \dots \ 1)$$

and $b = -D$. It is easy to see that for any \mathbf{e}_{ij} where $\mathbf{B}_{ij} = 1$, it holds that

$$\begin{aligned} \langle \mathbf{e}_{ij}, \mathbf{w} \rangle &= \alpha_i + (D - j + 1) \geq j + D - j + 1 = D + 1 \\ \implies H(\langle \mathbf{e}_{ij}, \mathbf{w} \rangle + b) &= H(\langle \mathbf{e}_{ij}, \mathbf{w} \rangle - D) = 1. \end{aligned}$$

On the other hand, for $\mathbf{B}_{ij} = 0$, we have

$$\begin{aligned} \langle \mathbf{e}_{ij}, \mathbf{w} \rangle &= \alpha_i + (D - j + 1) < j + D - j + 1 = D \\ \implies H(\langle \mathbf{e}_{ij}, \mathbf{w} \rangle + b) &= H(\langle \mathbf{e}_{ij}, \mathbf{w} \rangle - D) = 0. \end{aligned}$$

■

A more general useful class of detectable matrices are line matrices (Dewdney, 1977).

Definition 5.1.21. A binary matrix $\mathbf{B} \in \mathbb{B}^{D \times D}$ is a **line matrix** if any of the following conditions hold:

1. All \mathbf{B} 's ones lie either in the same row (\mathbf{B} is a **row matrix**) or in the same column (\mathbf{B} is a **column matrix**).
2. \mathbf{B} is a transversal, i.e., a matrix in which there is at most one 1 in any column and row.

Lemma 5.1.7. Row and column matrices are easily permutation-detectable.

Proof. Let $i, N \in \mathbb{Z}_D$ and \mathbf{B} be a row matrix with $\mathbf{B}_{ij_n} = 1$ for $n \in \mathbb{Z}_N$, i.e., a row matrix with all its ones in the i^{th} row. Define $\mathbf{P} \in \mathbb{B}^{D \times D}$ as $P_{1i} = 1$ and 0 elsewhere and $\mathbf{Q} \in \mathbb{B}^{D \times D}$ with $\mathbf{Q}_{j_nn} = 1$ and 0 elsewhere. Then, \mathbf{PBQ} contains all its 1 in its northwestern corner (contiguously in the first row) and is thus easily detectable. Let $\mathbf{w} \stackrel{\text{def}}{=} (\boldsymbol{\alpha} \mid D \ \dots \ 1)$, $b = D$ be the neuron weights from Lemma 5.1.6. Define $\mathbf{w}' \stackrel{\text{def}}{=} (\mathbf{P}^\top \boldsymbol{\alpha} \mid \mathbf{Q}(D \ \dots \ 1))$, $b' = D$. It is easy to see that this “rearranges” the components of the neuron recognizing the northwestern matrix \mathbf{PBQ} to make them recognize the original matrix, meaning that the neuron defined by \mathbf{w}' and b' recognizes the line matrix. The proof for a column matrix is analogous. ■

Lemma 5.1.8. Transversals are permutation-detectable.

Proof. The core idea of this proof is that every transversal can be permuted into a diagonal matrix, which can be written as a Hadamard product of a lower-triangular and an upper-triangular matrix.

Let \mathbf{B} be a transversal. Pre-multiplying \mathbf{B} with its transpose $\mathbf{P} \stackrel{\text{def}}{=} \mathbf{B}^\top$ results in a diagonal matrix. It is easy to see that \mathbf{PB} can be written as a Hadamard product $\mathbf{H}_1 \otimes \mathbf{H}_2$ of a lower-triangular matrix \mathbf{H}_1 and an upper-triangular matrix \mathbf{H}_2 . Both are easily permutation detectable. A conjunction of the neurons detecting \mathbf{H}_1 and \mathbf{H}_2 (again, performed by another neuron) detects the original matrix \mathbf{B} . In the following, we will refer to \mathbf{H}_1 and \mathbf{H}_2 as the *factors* of the transversal. ■

Crucially, any binary matrix $\mathbf{B} \in \mathbb{B}^{D \times D}$ can be decomposed into a set of line matrices \mathcal{B} whose disjunction is \mathbf{B} : $\bigvee_{\mathbf{M} \in \mathcal{B}} \mathbf{M} = \mathbf{B}$. It is easy to see that $\mathbf{B}_{ij} = 1$ if and only if there exists $\mathbf{M} \in \mathcal{B}$ such that $\mathbf{M}_{ij} = 1$. This means that non-zero entries of *any* $\mathbf{B} \in \mathbb{B}^{D \times D}$ decomposed into the set of line matrices \mathcal{B} can be detected using an MLP in two steps:

1. Detect the non-zero entries of the individual line matrices from the decomposition \mathcal{B} (which are, as shown above, detectable).
2. Take a disjunction of the detections of the individual line matrices to result in the activation of the original matrix.

The disjunction can again be performed by applying another 2-layer MLP to the activations of the line matrices. An important consideration in both Dewdney's as well as Indyk's construction later will be *how large* \mathcal{B} has to be.

Using matrix decomposition and detection for simulating the transition function. We now describe how Dewdney’s construction uses matrix detection based on the decomposition of matrices into line matrices to simulate an FSA using an HRNN. From a high level, the update steps of the HRNN will, just like in Minsky’s construction, simulate the transition function of the simulated FSA. However, in contrast to the Minsky construction, in which each transition step in the FSA was implemented by a *single* application of the Elman update rule, here, a single transition in the FSA will be implemented using *multiple* applications of the Elman update rule, the end result of which is the activation of the two-hot representation of the appropriate next state. Nonetheless, there are, abstractly, two sub-steps of the update step, analogous to the Minsky construction (cf. Fig. 5.6):

1. Detect the activations of all possible next states, considering any possible input symbol (performed by the term $\mathbf{U}\mathbf{h}_t$ in Minsky’s construction).
2. Filter the activations of the next states by choosing only the one transitioned into by a y_t -transition (performed by conjoining with the term $\mathbf{V}[[y_t]]$ in Minsky’s construction).

The novelty of Dewdney’s construction comes in the first sub-step: How can the Elman update step be used to activate the two-hot representation of q_t ’s out-neighborhood? As alluded to, this relies on the pre-computed predecessor matrices $\text{Pred}(q; y)$ (cf. Definition 5.1.16). The predecessor matrices of individual states are compressed (disjoined) into component-activating matrices, the representation matrices of the predecessors of specific *sets* of states (cf. Definition 5.1.17), defined through the function ϕ_2 in the following sense.

Definition 5.1.22. A **component-activating matrix** is the representation matrix $\mathbf{B}_{j,y,k} \stackrel{\text{def}}{=} \mathbf{W}_Q \text{Pred}(\Phi_{k,j}; y)$ for some $k \in \mathbb{Z}_r$ and $j \in \mathbb{Z}_2$.

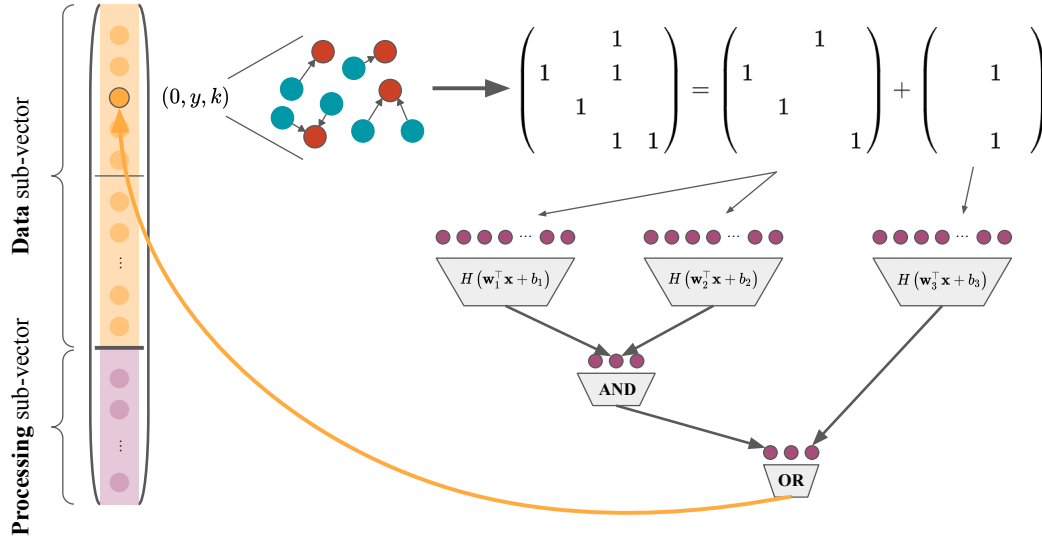
Intuitively, the component-activating matrix $\mathbf{B}_{j,y,k}$ is the result of the disjunction of the matrix representations of all y -predecessors q of all states q' whose j^{th} component of the vector $\phi_2(q')$ equals k . This results in $2|\Sigma|s$ matrices. They can be pre-computed and naturally depend on the transition function δ . The name component-activating matrix is inspired by the fact that each of the matrices “controls” the activation of one of the $2|\Sigma|s$ neurons in a specific sub-vector of the HRNN hidden state. That is, each component-activating matrix controls a particular dimension, indexed by the tuple (j, y, k) for $j \in \mathbb{B}, y \in \Sigma, k \in \mathbb{Z}_s$, in the **data sub-vector** of the HRNN hidden state. As we will see shortly, they contain all the information required for simulating \mathcal{A} with a HRNN.

To define the transition function of the HRNN simulating \mathcal{A} , all $2|\Sigma|s$ component-activating matrices are decomposed into permutation-detectable line matrices (cf. Definition 5.1.21) whose activations are combined (disjoined) into the activations of individual component-activating matrices. Analogously to above, we will denote the sets of line matrices decomposing the component-activating matrices as $\mathcal{B}_{j,y,k}$, i.e., $\mathbf{B}_{j,y,k} = \bigvee_{\mathbf{M} \in \mathcal{B}_{j,y,k}} \mathbf{M}$. The dimensions of the hidden state corresponding to the activations of the line matrices before they are combined into the activations of the component-activating matrices form the **processing sub-vector** of the HRNN hidden state since they are required in the pre-processing steps of the update step to determine the activation of the actual hidden state. This is schematically drawn in Fig. 5.12a.

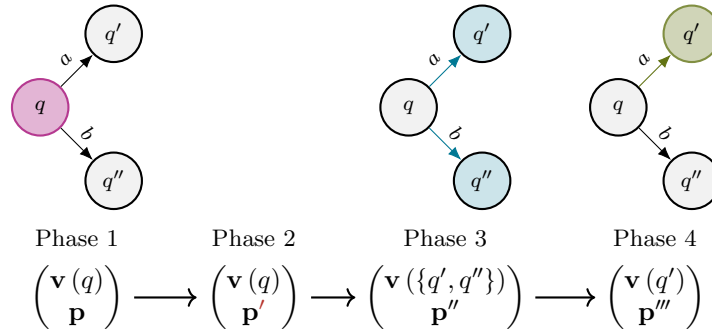
For any component-activating matrix \mathbf{B} decomposed into the set of line matrices \mathcal{B} , we know by Lemmas 5.1.7 and 5.1.8 that all $\mathbf{M} \in \mathcal{B}$ are detectable by a single-layer MLP. By adding an additional layer to the MLP, we can disjoin the detections of $\mathbf{M} \in \mathcal{B}$ into the detection of \mathbf{B} . More abstractly, this MLP, therefore, detects the activation of *one* of the $2|Q|s$ cells of the data sub-vector of the HRNN hidden state—all of them together then form the two-hot encoding of all possible next states of the FSA (before taking into account the input symbol). Designing $2|Q|s$ such single-values MLPs, therefore, results in an MLP activating the two-hot representations of all possible next states of the simulated FSA. Conjoining these activations with the input symbol, analogously to how this is done in the Minsky construction, results in the activation of the two-hot representation of only the actual next state of the simulated FSA. This is illustrated in Fig. 5.12b.

High-level overview of simulating a transition. In summary, after decomposing all the component-activating matrices into the sets $\mathcal{B}_{j,y,k}$, the detection of all candidate next states (before considering the input symbol) in the update step of HRNN is composed of the following sub-steps.

1. Compute the activations of the two *factors* of all the transversals in $\mathcal{B}_{j,y,k}$ for all j, y, k (Lemma 5.1.8).



(a) High-level overview of Dewdney's construction. The highlighted orange neuron in the representation of the state from the data sub-vector corresponds to the activation of one of the components of the red states (which have in common that their 0th component of $\phi_2(q)$ is the same). The matrix corresponding to the disjunction of the representations of their y -predecessors (blue states) is decomposed into two line matrices—a transversal and a column matrix. The non-zero elements of the former can be detected by a conjunction of two neurons while the non-zero elements of the latter can be detected directly by a single neuron. Those activations are then disjoined to result in the activation in the orange neuron. The purple neurons in the processing sub-vector are composed of the neurons in the networks implementing the detection of line matrices and their conjunctions and disjunctions (also shown in purple).



(b) A high-level illustration of how the transition function of the FSA is implemented in Dewdney's construction on an example of an FSA fragment, where the simulated automaton is initially in the state q and reads the symbol a , transitioning to q' . The components whose changes are relevant at a given step are highlighted. Starting in the state q , which is stored in the data sub-vector $\mathbf{v}(q)$, in the first sub-step, the processing bits of the appropriate line matrices are activated (\mathbf{p}'). Next, the activated line matrices are used to activate the representations of all the states in the out-neighborhood of q in the data sub-vector ($\mathbf{v}(\{q', q''\})$). Lastly, these representations are conjoined with the states reachable by the symbol a , resulting in the representation of the state q in the data sub-vector ($\mathbf{v}(q)$).

2. Conjoin the activations of the two factors into the activations of the transversals (Lemma 5.1.8).
3. Compute the activations of the column and row matrices in $\mathcal{B}_{j,y,k}$ for all j, y, k (Lemma 5.1.7).
4. Disjoin of the activations of all the line matrices (transversals, row, and column matrices) in $\mathcal{B}_{j,y,k}$ for all x, y, k to compute the activations of all $2|\Sigma|s$ component-activating matrices.

This results in the activation of the two-hot representations of all possible next states (i.e., the entire out-neighborhood of q_t). In the last sub-step of the HRNN update step, these are conjoined with the representation of the current input symbol. This step is very similar to the analogous stage in Minsky's construction, with the difference that here, the non-zero entries of the vector $\mathbf{V}\mathbf{h}_t$ must cover the two-hot representations of the states with an incoming y_t -transition. This conjunction then ensures that among all the states in the out-neighborhood of q_t , only the one reached by taking the y_t -transition will be encoded in \mathbf{h}_{t+1} . The construction just described can be summarized by the following lemma.¹⁸

Lemma 5.1.9. *Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be a deterministic FSA. Then, Dewdney's construction results in a HRNN correctly simulating \mathcal{A} 's transition function, i.e., $s(\mathbf{h}_t) = q_t$ for all t .*

This shows that Dewdney's construction correctly encodes the FSA in a HRNN. However, its space efficiency remains to be determined. As mentioned above, working with two-hot representations of the states means that the data sub-vector is of size $\mathcal{O}(|\Sigma|\sqrt{|Q|})$. However, the construction also requires a number of processing dimensions in the processing sub-vector. To understand the full complexity of the construction, we have to determine the maximal number of processing bits in the HRNN. The first step to the answer is contained in the following lemma, which describes the number of line matrices required to cover an arbitrary binary matrix. It lies in the core of the efficiency of Dewdney's construction.

Lemma 5.1.10. *Let $\mathbf{B} \in \mathbb{B}^{D \times D}$ with N^2 elements equalling 1. Then, there exists a decomposition \mathcal{B} of \mathbf{B} into at most $2N$ line matrices such that $\bigvee_{\mathbf{M} \in \mathcal{B}} \mathbf{M} = \mathbf{B}$.*

Proof. Based on Dewdney (1977). Define the sequence of transversals $\mathbf{T}_1, \mathbf{T}_2, \dots$ where \mathbf{T}_i is the transversal containing the maximum number of ones in the matrix $\mathbf{B}_i \stackrel{\text{def}}{=} \mathbf{B} - \bigvee_{j=1}^{i-1} \mathbf{B}_j$. The transversal containing the maximal number of ones can be found using the maximum matching algorithm. Continue this sequence until there are no more ones in \mathbf{B}_i . The number of ones in the matrices \mathbf{B}_i , $\|\mathbf{B}_i\|_1$, forms a (weakly) decreasing sequence.

If there are at most $2N$ transversals in the sequence, the lemma holds. Otherwise, we compare the functions $f(i) \stackrel{\text{def}}{=} \|\mathbf{T}_i\|_1$ and $g(i) \stackrel{\text{def}}{=} 2N - i$.

- If $f(i) > g(i)$ for all $i = 1, \dots, N$, then $\sum_{i=1}^N f(i) = \sum_{i=1}^N \|\mathbf{T}_i\|_1 > \sum_{i=1}^N 2N - i = 2N^2 - \frac{1}{2}N(N+1) \geq N^2$. However, the transversals in the decomposition cannot contain more ones than the original matrix.
- We conclude that for some $i \leq N$, $f(i) \leq g(i)$. Let i_0 be the first such index in $1, \dots, N$ and $\mathcal{L}_1 \stackrel{\text{def}}{=} \{\mathbf{T}_1, \dots, \mathbf{T}_{i_0}\}$. Since the maximum number of independent ones (in the sense that at most one appears in a single row/column) in \mathbf{B}_{i_0-1} is $\|\mathbf{T}_{i_0}\|_1 \leq 2N - i_0$ (those are chosen by the maximum transversal \mathbf{T}_{i_0}). By König's theorem (Szárnyas, 2020), there is a set of at most $2N - i_0$ column or row matrices $\mathcal{L}_2 \stackrel{\text{def}}{=} \{\mathbf{L}_1, \dots, \mathbf{L}_k\}$ with $k \leq 2N - i_0$ which cover \mathbf{B}_{i_0-1} .¹⁹ Therefore, $\mathcal{L} \stackrel{\text{def}}{=} \mathcal{L}_1 \cup \mathcal{L}_2$ constitutes a valid cover of \mathbf{B} with $\leq N + 2N - i_0 = \mathcal{O}(N)$ matrices. ■

We will denote the number of matrices in the line decomposition of a matrix \mathbf{B} constructed by the greedy procedure from Lemma 5.1.10 as $L(\mathbf{B})$. Connecting this lemma to Dewdney's construction, this shows that the number of neurons required to detect the activation of a *single set* $\text{Pred}(k; y)$ grows asymptotically as the square root of the number of ones in the representation matrix $\mathbf{W}_Q \text{Pred}(k; y)$ —this is how many line matrices the matrix will decompose into. The size of each neuron is $2|\Sigma|s$.

¹⁸To formally prove it is correct, we would have to follow a similar set of steps to how the correctness of Minsky's construction (Lemma 5.1.4) was proved. We omit this for conciseness.

¹⁹Intuitively, since all ones are contained within $\leq 2N - i_0$ rows or columns, they can be simply covered by matrices containing those.

This allows us to show how many neurons the entire HRNN simulating \mathcal{A} has. Since we know that the data sub-vector will always have exactly $2|\Sigma|s$ cells, we characterize the number of processing cells in the following lemma.

Lemma 5.1.11. *Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be a deterministic FSA. Then, Dewdney's construction results in a HRNN with a hidden state of size $\mathcal{O}\left(|\Sigma||Q|^{\frac{3}{4}}\right)$.*

Proof. The number of cells in the entire processing sub-vector is simply the sum of the processing neurons of all the data components. In the worst case, a single component-activating matrix \mathbf{B} requires $2L(\mathbf{B}) + 1$ neurons (2 for each transversal in the decomposition of \mathbf{B} and an additional one for their disjunction). Therefore, enumerating the set of matrices $\{\mathbf{B}_{j,y,k} \mid j \in \mathbb{Z}_2, y \in \Sigma, k \in \mathbb{Z}_s\}$ with \mathbf{B}_n for $n = 1, \dots, 2|\Sigma|s$, the number of neurons required by all component-activating matrices is bounded as follows.

$$\sum_{n=1}^{2|\Sigma|s} 2L(\mathbf{B}_n) + 1 \leq \sum_{n=1}^{2|\Sigma|s} 2\left(2\lceil\sqrt{\|\mathbf{B}_n\|_1}\rceil\right) + 1 \stackrel{\text{def}}{=} \sum_{n=1}^{2|\Sigma|s} 4m_n + 1 \quad (5.82)$$

Since the matrices \mathbf{B}_n contain one non-zero entry for each state-symbol pair, it holds that

$$\sum_{n=1}^{2|\Sigma|s} \|\mathbf{B}_n\|_1 \leq \sum_{n=1}^{2|\Sigma|s} m_n^2 = |\Sigma||Q| \quad (5.83)$$

Pretending that m_n can take real values, the value of Eq. (5.82) is maximized under the constraint from Eq. (5.83) when all m_n are equal with $m_n = \sqrt{2s}$. This means that

$$\sum_{n=1}^{2|\Sigma|s} 4m_n + 1 \leq \sum_{n=1}^{2|\Sigma|s} 4\sqrt{2s} + 1 = 8|\Sigma|s\sqrt{2s} + 1 = \mathcal{O}\left(|\Sigma||Q|^{\frac{3}{4}}\right), \quad (5.84)$$

finishing the proof. ■

All results stated in this section can be summarized in the following theorem.

Theorem 5.1.5. *Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be a deterministic FSA. Then, there exists a HRNN of size $\mathcal{O}\left(|\Sigma||Q|^{\frac{3}{4}}\right)$ correctly simulating \mathcal{A} .*

Indyk's Construction

§5.1.7 describes a construction of an HRNN of size $\mathcal{O}\left(|\Sigma||Q|^{\frac{3}{4}}\right)$ simulating an FSA. While this improves the space efficiency compared to Minsky's construction, it is not asymptotically optimal. Indyk (1995) proved that a HRNN simulating an FSA $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ over a binary alphabet $\Sigma = \mathbb{B}$ requires at least $\Omega\left(\sqrt{|Q|}\right)$ hidden dimensions. He also provided a construction that achieves this lower bound. This construction is conceptually very similar to Dewdney's in that it works by activating neurons corresponding to some form of compressed predecessor matrices (component-activating matrices) and then selecting the transition which matches the input symbol. Again, it additively covers these matrices with components that are easy to detect, similar to how Dewdney's construction uses line matrices. However, Indyk's construction defines component-activating matrices based on different sets of states and covers them with a different decomposition—these are the two crucial differences allowing the construction to achieve the optimal lower bound.

We first define the component-activating matrices and their role in updating the hidden state of the HRNN. In Indyk's construction, the component-activating matrices are based on *four-hot* rather than two-hot encodings of states.

Definition 5.1.23. *Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be an FSA, $r \stackrel{\text{def}}{=} \lceil |Q|^{\frac{1}{4}} \rceil$, and π a permutation of $Q = [|Q|]$.²⁰ We define the **four-hot representation** of $q \in Q$ as*

$$\phi_4(q) = (\ell_1, \ell_2, \ell_3, \ell_4) \quad (5.85)$$

²⁰The exact form of π will be important later. For now, one can think of π as the identity function.

where

$$\ell_j = \frac{\pi(q)}{r^{j-1}} \pmod{r}. \quad (5.86)$$

We denote the inverse of ϕ_4 with ϕ_4^{-1} and further define for $k \in \mathbb{Z}_r$

$$\phi_4^{-1}(k, \cdot, \cdot, \cdot) \stackrel{\text{def}}{=} \{q \in Q \mid \phi_4(q)_1 = k\} \quad (5.87)$$

and $\phi_4^{-1}(\cdot, k, \cdot, \cdot)$, $\phi_4^{-1}(\cdot, \cdot, k, \cdot)$, and $\phi_4^{-1}(\cdot, \cdot, \cdot, k)$ analogously.

We will denote $\phi_4^{-1}(\dots, k, \dots)$ with k in j^{th} position (with $j \in \mathbb{Z}_4$) as $\Phi_{k,j}$. Despite using the four-hot representations, Indyk's construction still requires the two-hot representations based on ϕ_2 as before. In this case, however, they again depend on the chosen permutation π . This allows us to define the component-activating matrices as follows.

Definition 5.1.24. A **component-activating matrix** in Indyk's construction is the representation matrix $\mathbf{W}_{Q\text{Pred}}(\Phi_{k,j}; y)$ for some $k \in \mathbb{Z}_r$, $j \in \mathbb{Z}_4$, and $y \in \Sigma$.

For efficient detection, the component-activating matrices are covered by so-called non-decreasing matrices.

Definition 5.1.25. We say that $\mathbf{B} \in \mathbb{B}^{D \times D}$ is **non-decreasing** if there exists a non-decreasing (partial) function $f: \mathbb{Z}_D \rightarrow \mathbb{Z}_D$ (from columns to rows) such that

$$\mathbf{B}_{ij} = 1 \iff f(j) = i \quad (5.88)$$

and, if f is defined for some $j \in \mathbb{Z}_D$, it is also defined for all $j' \geq j$.

Example 5.1.15. An example of a non-decreasing matrix is

$$\mathbf{B} = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}. \quad (5.89)$$

The function f defining the non-decreasing matrix \mathbf{B} is $f = \begin{pmatrix} 0 & 1 & 2 & 3 \\ \emptyset & 0 & 1 & 1 \end{pmatrix}$, where \emptyset denotes that the function is not defined.

Again, clearly, any matrix $\mathbf{B} \in \mathbb{B}^{D \times D}$ can be (non-uniquely) decomposed into at most D non-decreasing matrices. Moreover, non-decreasing matrices are detectable.

Lemma 5.1.12. Non-decreasing matrices are detectable.

Proof. Let $\mathbf{B} \in \mathbb{B}^{D \times D}$ be a non-decreasing matrix defined by the partial function f . Divide the domain of f into the set of intervals in which the function is constant, with $I(j)$ denoting the interval of $j \in \mathbb{Z}_{r^2}$ for j such that $f(j)$ is defined. Then, it is easy to see that $\mathbf{B}_{ij} = 1 \iff i = f(j)$, meaning that by defining the parameters \mathbf{w} and b as

$$w_{f(j)} \stackrel{\text{def}}{=} r^2 - I(j) \quad (5.90)$$

$$w_{r^2+j} \stackrel{\text{def}}{=} I(j) \quad (5.91)$$

$$b \stackrel{\text{def}}{=} -r^2 \quad (5.92)$$

and other elements as 0, we get that

$$\mathbf{B}_{ij} = 1 \iff i = f(j) \iff w_i + w_j + b = 0. \quad (5.93)$$

Compared to earlier, where component-activating matrices were detected by testing an *inequality*, detecting a non-decreasing matrix requires testing an *equality*. Since all terms in the equality are integers, testing the equality can be performed with the Heaviside activation function by conjoining two neurons; one testing the inequality $w_i + w_j + b - 1 < 0$ and another one testing the inequality $w_i + w_j + b + 1 > 0$. Both can individually be performed by a single neuron and then conjoined by an additional one. ■

With this, the high-level idea of Indyk’s construction is outlined in Fig. 5.13. After constructing the component-activating matrices based on ϕ_4 and decomposing them into non-decreasing matrices, the rest of Indyk’s construction is very similar to Dewdney’s construction, although the full update step of the HRNN requires some additional processing. To test the equality needed to detect non-decreasing matrices in the decomposition, Eq. (5.93), the four-hot representations are first converted into two-hot ones. This can be done by a simple conjunction of the first two and the last two components of the four-hot representation. Then, the activations of the non-decreasing matrices can be computed and disjointed into the representations of the component-activating matrices. These form the $4|\Sigma|r$ components of the data sub-vector of the HRNN hidden state. They contain the activations of all possible next states, i.e., the out-neighborhood of the current state of \mathcal{A} . These are then conjoined with the representation of the current input symbol in the same way as in Dewdney’s construction but adapted to the four-hot representations of the states. The process is thus very similar to the phases of Dewdney’s construction illustrated in Fig. 5.12b.

Indyk’s construction can be summarized by the following lemma.²¹

Lemma 5.1.13. *Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be a deterministic FSA. Then, Indyk’s construction results in a HRNN correctly simulating \mathcal{A} ’s transition function, i.e., $s(\mathbf{h}_t) = q_t$ for all t .*

The only remaining thing to show is that Indyk’s construction achieves the theoretically optimal lower bound on the size of the HRNN simulating a deterministic FSA. All previous steps of the construction were valid no matter the chosen permutation π . The permutation, however, matters for space efficiency: intuitively, it determines how efficiently one can decompose the resulting component-activating matrices (which depend on the permutation) into non-decreasing matrices in the sense of how many non-decreasing matrices are required to cover it. Indyk, therefore, proved that there always exists, with non-zero probability, a permutation in which the decomposition across all states is efficient enough to achieve the minimum number of neurons required. This is formalized by the following lemma, whose proof can be found in Indyk (1995, Lemma 6).

Lemma 5.1.14. *Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be a deterministic FSA. There exists a permutation of Q such that Indyk’s construction results in a HRNN of size $\mathcal{O}(|\Sigma|\sqrt{|Q|})$.*

This concludes our presentation of Indyk’s construction. All results stated in this section can be summarized by the following theorem.

Theorem 5.1.6. *Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be a deterministic FSA. There exists a HRNN of size $\mathcal{O}(|\Sigma|\sqrt{|Q|})$ correctly simulating \mathcal{A} .*

Lower Bound in the Probabilistic Setting

We now ask whether the same lower bound can also be achieved when simulating dPFSA. We find that the answer is negative: dPFSA may require an HRNN LMs of size $\Omega(|\Sigma||Q|)$ to faithfully represent their probability distribution. Since the transition function of the underlying FSA can be simulated with more efficient constructions, the bottleneck comes from defining the same probability distribution. Indeed, as the proof of the following theorem shows, the issue intuitively arises in the fact that, unlike in an HRNN LM, the local probability distributions of the different states in a PFSA are completely arbitrary, whereas they are defined by shared parameters (the output matrix \mathbf{E}) in an HRNN LM.

Theorem 5.1.7. *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a minimal dPFSA and $(\Sigma, D, \mathbf{U}, \mathbf{V}, \mathbf{E}, \mathbf{b}_h, \mathbf{h}_0)$ an HRNN LM defining the same LM. Then, D must scale linearly with $|Q|$.*

Proof. Without loss of generality, we work with \mathbb{R} -valued hidden states. Let \mathcal{A} be a minimal deterministic PFSA and $\mathcal{R} = (\Sigma, D, \mathbf{U}, \mathbf{V}, \mathbf{E}, \mathbf{b}_h, \mathbf{h}_0)$ a HRNN with $p_{\text{LN}}(\mathbf{y}) = \mathcal{A}(\mathbf{y})$ for every $\mathbf{y} \in \Sigma^*$. Let $\mathbf{y}_{<T} \in \Sigma^*$ and $\mathbf{y}_{\leq T} \stackrel{\text{def}}{=} \mathbf{y}_{<T}y$ for some $y \in \Sigma$. Define $\bar{p}(\mathbf{y}) \stackrel{\text{def}}{=} \prod_{t=1}^{|\mathbf{y}|} p_{\text{SM}}(y_t | \mathbf{y}_{<t})$. It is easy to see that $\bar{p}(\mathbf{y}_{<T}y_T) = \bar{p}(\mathbf{y}_{<T})p_{\text{SM}}(y_T | \mathbf{y}_{<T})$. The conditional distribution $p_{\text{SM}}(\cdot | \mathbf{y}_{<T})$ are proportional to the values in $\mathbf{E}\mathbf{h}_{T-1}$.

²¹Again, to formally prove it is correct, we would have to follow a similar set of steps to how the correctness of Minsky’s construction (Lemma 5.1.4) was proved. We omit this for conciseness.

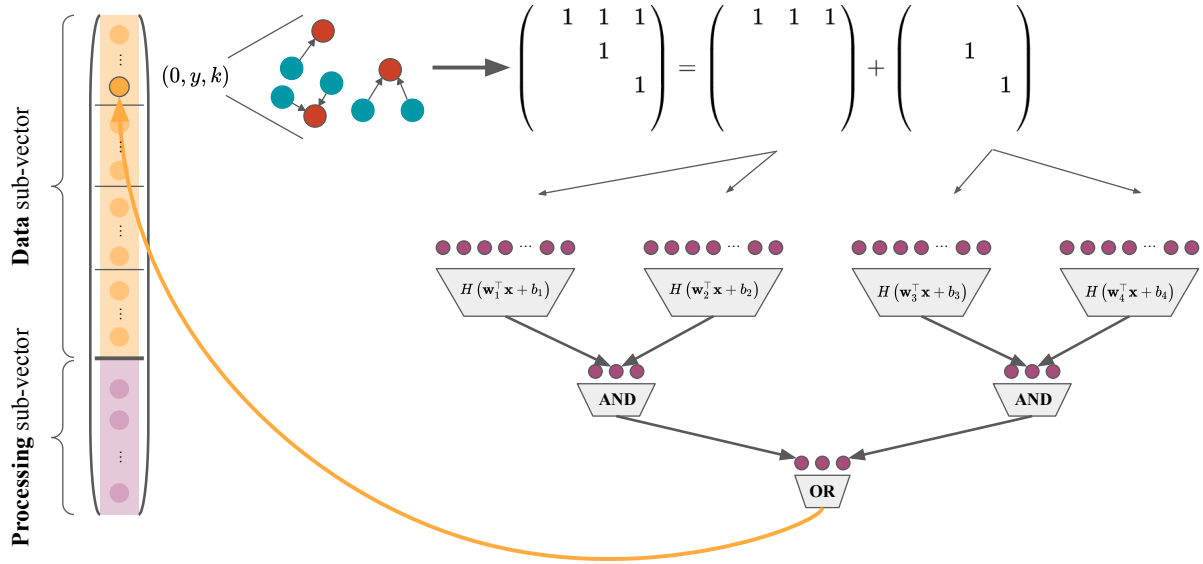
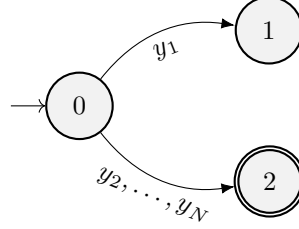


Figure 5.13: High-level overview of Indyk's construction. The highlighted **orange** neuron in the representation of the state from the data sub-vector corresponds to the activation of one of the components of the **red** states (which have in common that their 0^{th} component of $\phi_4(q)$ is the same). The matrix corresponding to the disjunction of the representations of their y -predecessors (**blue** states) is decomposed into two non-decreasing matrices. The non-zero elements of both can be detected by a conjunction of two neurons; here, $f_1 = \begin{pmatrix} 0 & 1 & 2 & 3 \\ \emptyset & 0 & 0 & 0 \end{pmatrix}$ and $f_2 = \begin{pmatrix} 0 & 1 & 2 & 3 \\ \emptyset & \emptyset & 1 & 2 \end{pmatrix}$, meaning that $\mathbf{w}_1 = (3 \ 0 \ 0 \ 0 \mid 0 \ 1 \ 1 \ 1)$, $\mathbf{w}_2 = (0 \ 3 \ 2 \ 0 \mid 0 \ 0 \ 1 \ 2)$, and $b_1 = b_2 = 4$. Those activations are then disjoined to result in the activation in the **orange** neuron. The **purple** neurons in the processing sub-vector are composed of the neurons in the networks implementing the detection of line matrices and their conjunctions and disjunctions (also shown in **purple**). Note that even if the second matrix were not non-decreasing in itself (i.e., the columns of the two ones would be flipped), one could still transform it into a non-decreasing matrix by permuting the columns and permuting the corresponding neurons.

Figure 5.14: The FSA \mathcal{A}_N .

By definition of the deterministic PFSA, there are $|Q|$ such conditional distributions. Moreover, these distributions (represented by vectors $\in \Delta^{|\Sigma|-1}$) can generally be *linearly independent*. This means that for any q , the probability distribution of the outgoing transitions can not be expressed as a linear combination of the probability distributions of other states. To express the probability vectors for all states, the columns of the output matrix \mathbf{E} , therefore, have to span $\mathbb{R}^{|Q|}$, implying that \mathbf{E} must have at least $|Q|$ columns. This means that the total space complexity (and thus the size of the HRNN representing the same distribution as \mathcal{A}) is $\Omega(|Q|)$. ■

Note that there exist regular languages motivated by phenomena in human language that *can* be represented in logarithmic space in the number of the states of their minimal FSA. For example, [Hewitt et al. \(2020\)](#) show that bounded Dyck languages of k parenthesis types and of maximal depth m , which require an FSA with k^m states to be recognized, can be represented by HRNN LMs of size $m \log k$, which is an exponential improvement on Indyk’s lower bound.

Asymptotic Bounds in $|\Sigma|$ Since each of the input symbols can be encoded in $\log |\Sigma|$ bits, one could expect that the linear factor in the size of the alphabet from the constructions above could be reduced to $\mathcal{O}(\log |\Sigma|)$. However, we again find that such reduction is in general not possible—the set of FSAs presented next is an example of a family that requires an HRNN whose size scales linearly with $|\Sigma|$ to be simulated correctly. We also provide a sketch of the proof of why a compression in $|\Sigma|$ is not possible.

Let $\mathcal{A}_N = (\Sigma_N, \{0, 1\}, \{0\}, \{1\}, \delta_N)$ be an FSA over the alphabet $\Sigma_N = \{y_1, \dots, y_N\}$ such that $\delta_N = \{0 \xrightarrow{y_1} 1\} \cup \{0 \xrightarrow{y_n} 2 \mid n = 2, \dots, N\}$ (see Fig. 5.14).

Clearly, to be able to correctly represent all local distributions of the dPFSA, the HRNN LM must contain a representation of each possible state of the dPFSA in a unique hidden state. On the other hand, the only way that the HRNN can take into account the information about the current state q_t of the simulated FSA \mathcal{A} is through the hidden state \mathbf{h}_t . The hidden state, in turn, only interacts with the recurrence matrix \mathbf{U} , which does not have access to the current input symbol y_{t+1} . The only interaction between the current state and the input symbol is thus through the addition in $\mathbf{U}\mathbf{h}_t + \mathbf{V}\llbracket y_{t+1} \rrbracket$. This means that, no matter how the information about q_t is encoded in \mathbf{h}_t , in order to be able to take into account all possible transitions stemming in q_t (before taking into account y_{t+1}), $\mathbf{U}\mathbf{h}_t$ must activate *all* possible next states, i.e., the entire out-neighborhood of q_t . On the other hand, since $\mathbf{V}\llbracket y_{t+1} \rrbracket$ does not have precise information about q_t , it must activate all states which can be entered with an y_{t+1} -transition, just like in Minsky’s construction.

In Minsky’s construction, the recognition of the correct next state was done by keeping a separate entry (one-dimensional sub-vector) for each possible pair q_{t+1}, y_{t+1} . However, when working with compressed representations of states (e.g., in logarithmic space), a single common sub-vector of size $< |\Sigma|$ (e.g., $\log |\Sigma|$) has to be used for all possible symbols $y \in \Sigma$. Nonetheless, the interaction between $\mathbf{U}\mathbf{h}_t$ and $\mathbf{V}\llbracket y_{t+1} \rrbracket$ must then ensure that only the correct state q_{t+1} is activated. For example, in Minsky’s construction, this was done by simply taking the conjunction between the entries corresponding to q, y in $\mathbf{U}\mathbf{h}_t$ and the entries corresponding to q', y' in $\mathbf{V}\llbracket y' \rrbracket$, which were all represented in individual entries of the vectors. On the other hand, in the case of the log encoding, this could intuitively be done by trying to match the $\log |\Sigma|$ ones in the representation $(\mathbf{p}(y) \mid \mathbf{1} - \mathbf{p}(y))$, where $\mathbf{p}(y)$ represent the binary encoding of y . If the $\log |\Sigma|$ ones match (which is checked simply as it would result in a large enough sum in the corresponding entry of the matrix-vector product), the correct transition could be chosen (to perform the conjunction from Fact 5.1.4

correctly, the bias would simply be set to $\log |\Sigma| - 1$). However, an issue arises as soon as *multiple* dense representations of symbols in $\mathbf{V}[[y]]$ have to be activated against the same sub-vector in $\mathbf{U}\mathbf{h}_t$ —the only way this can be achieved is if the sub-vector in $\mathbf{U}\mathbf{h}_t$ contains the disjunction of the representations of all the symbols which should be activated with it. If this sets too many entries in $\mathbf{U}\mathbf{h}_t$ to one, this can result in “false positives”. This is explained in more detail for the dPFSA in Fig. 5.14 next.

Let \mathbf{r}_n represent any dense encoding of y_n in the alphabet of \mathcal{A}_N (e.g., in the logarithmic case, that would be $(\mathbf{p}(n) \mid \mathbf{1} - \mathbf{p}(n))$). Going from the intuition outlined above, any HRNN simulating \mathcal{A}_N , the vector $\mathbf{U}\mathbf{h}_0$ must, among other things, contain a sub-vector corresponding to the states 1 and 2. The sub-vector corresponding to the state 2 must activate (through the interaction in the Heaviside function) against any y_n for $n = 2, \dots, N$ in \mathcal{A}_N . This means it has to match all representations \mathbf{r}_n for all $n = 2, \dots, N$. The only way this can be done is if the pattern for recognizing state 2 being entered with any y_n for $n = 2, \dots, N$ is of the form $\mathbf{r} = \bigvee_{n=2}^N \mathbf{r}_n$. However, for sufficiently large N , $\mathbf{r} = \bigvee_{n=2}^N \mathbf{r}_n$ will be a vector of all ones—including all entries active in \mathbf{r}_1 . This means that *any* encoding of a symbol will be activated against it—among others, y_1 . Upon reading y_1 in state 1, the network will therefore not be able to deterministically activate only the sub-vector corresponding to the correct state 1. This means that the linear-size encoding of the symbols is, in general, optimal for representing dPFSA with HRNN LMs. This discussion implies the following theorem.

Theorem 5.1.8. *Let $\mathcal{A} = (\Sigma, Q, \delta, \lambda, \rho)$ be a minimal dPFSA and $(\Sigma, D, \mathbf{U}, \mathbf{V}, \mathbf{E}, \mathbf{b}_h, \mathbf{h}_0)$ an HRNN LM defining the same LM. Then, D must scale linearly with $|\Sigma|$.*

Based on the challenges encountered in the example above, we can devise a simple sufficient condition for a logarithmic compression w.r.t. $|\Sigma|$ to be possible: namely, that for any pair of states $q, q' \in Q$, there is at most a single transition leading from q to q' . Importantly, this condition is met by classical n -gram LMs, and also in the language studied by Hewitt et al. (2020), which allowed them to reduce the complexity w.r.t. $|\Sigma|$ to the logarithmic factor. This intuitive characterization can be formalized by a property we call $\log |\Sigma|$ -separability.

Definition 5.1.26. *An FSA $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ is $\log |\Sigma|$ -separable if it is deterministic and, for any pair $q, q' \in Q$, there is at most one symbol $y \in \Sigma$ such that $q \xrightarrow{y} q' \in \delta$.*

$\log |\Sigma|$ -separability is a relatively restrictive condition. To amend that, we introduce a simple procedure which, at the expense of enlarging the state space by a factor of $|\Sigma|$, transforms a general deterministic (unweighted) FSA into a $\log |\Sigma|$ -separable one. We call this $\log |\Sigma|$ -separation. Intuitively, it augments the state space by introducing a new state (q, y) for every outgoing transition $q \xrightarrow{y} q'$ of every state $q \in Q$, such that (q, y) simulates the only state the original state q would transition to upon reading y . Due to the determinism of the original FSA, this results in a $\log |\Sigma|$ -separable FSA with at most $|Q||\Sigma|$ states.

While the increase of the state space might seem like a step backward, recall that using Indyk’s construction, we can construct an HRNN simulating an FSA whose size scales with the square root of the number of states. And, since the resulting FSA is $\log |\Sigma|$ -separable, we can reduce the space complexity with respect to Σ to $\log |\Sigma|$. This is summarized in the following theorem, which characterizes how compactly general deterministic FSAs can be encoded by HRNNs. To our knowledge, this is the tightest bound on simulating general unweighted deterministic FSAs with HRNNs.

Theorem 5.1.9. *Let $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ be a minimal FSA recognizing the language L . Then, there exists an HRNN $\mathcal{R} = (\Sigma, D, \mathbf{U}, \mathbf{V}, \mathbf{E}, \mathbf{b}_h, \mathbf{h}_0)$ accepting L with $D = \mathcal{O}\left(\log |\Sigma| \sqrt{|\Sigma||Q|}\right)$.*

The full $\log |\Sigma|$ -separation procedure is presented in Algorithm 2. It follows the intuition of creating a separate “target” for each transition $q \xrightarrow{y} q'$ for every state $q \in Q$. To keep the resulting FSA deterministic, a new, artificial, initial state with no incoming transitions is added and is connected with the augmented with the out-neighborhood of the original initial state.

The following simple lemmata show the formal correctness of the procedure and show that it results in a $\log |\Sigma|$ -separable FSA, which we need for compression in the size of the alphabet.

Lemma 5.1.15. *For any $y \in \Sigma$, $(q, y) \xrightarrow{y'} (q', y') \in \delta'$ if and only if $q \xrightarrow{y'} q' \in \delta$.*

Proof. Ensured by the loop on Line 3. ■

Algorithm 2

```

1. def Separate( $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ ):
2.    $\mathcal{A}' \leftarrow (\Sigma, Q' = Q \times \Sigma \cup \{q_\epsilon'\}, \delta' = \emptyset, I' = \{q_\epsilon'\}, F' = \emptyset)$ 
3.    $\triangleright$  Connect the out-neighborhood of the original initial state  $q_\epsilon$  with the new, artificial, initial state.
4.   for  $y \in \Sigma$  :
5.     for  $q_\epsilon \xrightarrow{y'} q' \in \delta$  :
6.       add  $q_\epsilon' \xrightarrow{y} (q', y')$  to  $\delta'$ 
7.   for  $q \in Q, y \in \Sigma$  :
8.     for  $q \xrightarrow{y'} q' \in \delta$  :
9.       add  $(q, y) \xrightarrow{y'} (q', y')$  to  $\delta'$ 
10.   $\triangleright$  Add all state-symbol pairs with a state from the original set of final states to the new set of final states.
11.  for  $q_\varphi \in F, y \in \Sigma$  :
12.    add  $(q_\varphi, y)$  to  $F'$ 
13.  if  $q_\epsilon \in I$  :  $\triangleright$  Corner case: if the original initial state  $q_\epsilon$  is an initial state, make the artificial initial state  $q_\epsilon'$  final.
14.    add  $q_\epsilon'$  to  $F'$ 
15.  return  $\mathcal{A}'$ 

```

Lemma 5.1.16. $\log |\Sigma|$ -separation results in an equivalent FSA.

Proof. We have to show that, for any $\mathbf{y} \in \Sigma^*$, \mathbf{y} leads to a final state in \mathcal{A} if and only if \mathbf{y} leads to a final state in \mathcal{A}' . For the string of length 0, this is clear by Lines 13 and 14. For strings of length ≥ 1 , it follows from Lemma 5.1.15 that \mathbf{y} leads to a state q in \mathcal{A} if and only if $\exists y \in \Sigma$ such that \mathbf{y} leads to (q, y) in \mathcal{A}' . From Lines 11 and 12, $(q, y) \in F'$ if and only if $q \in F$, finishing the proof. ■

Lemma 5.1.17. $\log |\Sigma|$ -separation results in a $\log |\Sigma|$ -separable FSA.

Proof. Since the state (q', y') is the only state in Q' transitioned to from (q, y) after reading y' (for any $y \in \Sigma$), it is easy to see that \mathcal{A}' is indeed $\log |\Sigma|$ -separable. ■

Discussion and the practical applicability of these result. This section showed that Heaviside-activated RNNs are equivalent to WFSAs. This might come as a bit of a surprise considering that we introduced RNNs with the goal of overcoming some limitations of exactly those models, e.g., the finite context length. However, note that to arrive at this result, we considerably restricted the form of a recurrent neural network. While on the one hand restriction to the Heaviside activation function means that all the RNNs we considered in this section can be implemented and represented in a computer, the RNN sequence models that we usually deal with are much more complex than this analysis allowed for. Furthermore, note that the RNNs in practice do not learn sparse hidden states of the form considered in the construction in the proof of Lemma 5.1.3—indeed, networks with Heaviside activation functions are not trainable with methods discussed in §3.2 as the gradient on the entire parameter space would be either 0 or undefined and in this sense, the trained networks would never have such hidden state dynamics. The dynamics of RNNs in practice result in *dense* hidden states, i.e., states in which many dimensions are non-zero. Nonetheless, keep in mind that theoretically, due to the finite-precision nature of our computers, all models we ever consider will be at most finite-state—the differentiating factor between them will be how appropriately to the task they are able to learn the topology (transitions) of the finite-state automaton they represent and how efficiently they are able to learn it.

Turing Completeness of Recurrent Neural Networks

We now turn to the (purely theoretical) treatment of the expressive capacity of recurrent neural networks in which we take the liberty of making somewhat unrealistic assumptions. Specifically, in practice, RNNs usually have the following properties:

- The weights and intermediate computations are done with finite floating point precision;
- An RNN always operates in *real-time*, meaning that it performs a constant number of operations before consuming/outputting a symbol.

Under these assumptions, we saw that RNNs with Heaviside activations in a practical setting lie at the bottom of the weighted Chomsky hierarchy, being able to only recognize regular languages. However, if we relax these two assumptions, allowing for arbitrary precision and unbounded computation time between symbols, RNNs jump directly to the top of the hierarchy: they become Turing complete.

We start by introducing the saturated sigmoid, one of the building blocks we will use to show this.

Definition 5.1.27 (Saturated Sigmoid). *The **saturated sigmoid** is defined as*

$$\sigma(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } 0 < x \leq 1 \\ 1 & \text{if } x > 1 \end{cases} \quad (5.94)$$

Intuitively, the saturated sigmoid clips all negative values to 0, all values larger than 1 to 1, and leaves the elements of $[0, 1]$ intact. The graph of this function is shown in Fig. 5.15.

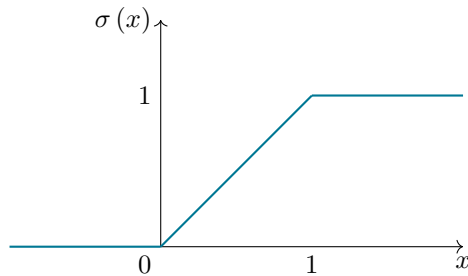


Figure 5.15: The saturated sigmoid.

The central result of this subsection is then summarized in the following theorem.

Theorem 5.1.10 (Saturated Sigmoid Elman RNNs are Turing complete). *Elman recurrent neural network sequence models with the saturated sigmoid activation functions are Turing complete.*

By the end of this subsection, we will have proven this result by showing that Saturated Sigmoid Elman RNNs can encode two-stack pushdown automata which are computationally equivalent to Turing machines (cf. Definitions 4.2.50 and 4.2.51). We start with a simpler construction: building on our placement of RNNs on at least the regular rung of the ladder of formal language complexity (cf. Lemma 5.1.3), we take one step up the ladder and show that RNNs can simulate a *single-stack* pushdown automaton (cf. Definition 4.2.31). This will help us gain the intuition behind how RNNs can use infinite precision arithmetic to simulate a stack. We will then simply generalize this construction to the two-stack case.

Let us begin by considering the problem of representing a stack—an arbitrarily long sequence of symbols, accessible in a last-in-first-out (LIFO) fashion—in a vector of constants size, e.g., the hidden state of a recurrent neural network. For simplicity, but without the loss of generality, assume that we are working with a simple two-letter stack alphabet $\Gamma = \{0, 1\}$. Any stack sequence γ will be a member of Γ^* , i.e., a string of 0's and 1's. If we think of the stack *symbols* as *numbers* for a moment, there is a natural correspondence between the possible stack configurations and *numbers* expressed in base 2. By convention, we will represent a string of stack symbols γ with numbers after the decimal point, rather than as integers. Assuming infinite

precision, we can therefore simply represent each stack configuration as a single number (of course, the stack alphabet does not have to be exactly $\Gamma = \{0, 1\}$ —we can always map symbols from any alphabet into their numeric representations in some base large enough to allow for the entire alphabet). Notice that in this case, pushing or popping from the stack can be performed by division and multiplication of the value representing the stack—if we want to push a value $x \in \{0, 1\}$, we can *divide* the current representation (by 2) and append x to the right side of the new representation and if we want to pop any value, we simply have to *multiply* the current representation by 2. This also gives us an idea of how to represent a stack in the hidden state of an RNN: the *entire* stack sequence will simply be represented in a single dimension of the hidden state, and the value stored in the cell will be updated according to the transitions defined by the simulated automaton. Note, however, that the RNN will not only have a single dimension in the hidden state: other dimensions will contain values that will be required to control the RNN updates correctly.

In our proofs, we consider a special type of pushdown automata, as defined in Definition 4.2.31: we will use pushdown automata which only consider the *topmost* element of the stack when defining the possible transitions from a configuration and can only push one stack symbol at a time. More formally, this means that in the tuple $\mathcal{P} = (\Sigma, Q, \Gamma, \delta, (q_\iota, \gamma_\iota), (q_\varphi, \gamma_\varphi))$, we have that $\delta \subseteq Q \times \Gamma \times (\Sigma \cup \{\varepsilon\}) \times Q \times \Gamma$ rather than the more general $\delta \subseteq Q \times \Gamma^* \times (\Sigma \cup \{\varepsilon\}) \times Q \times \Gamma^*$. Furthermore, we assume that $\gamma_\iota = \varepsilon$ and $\gamma_\varphi = \varepsilon$, that is, the PDA starts off with an empty stack and has to empty it again to arrive at a final configuration. Note that these restrictions can be done without loss of generality—that is, such pushdown automata are as powerful as the unrestricted versions (Sipser, 2013). With this in mind, we can show that arbitrary precision RNNs are capable of recognizing at least deterministic context-free languages:

Theorem 5.1.11. *Elman recurrent neural networks can recognize deterministic context-free languages.*

Before we continue to the proof of Theorem 5.1.11, let us remark on three simple but important intuitions which will be crucial for understanding the construction of the Elman RNN, both in the single- as well as the two-stack variants of PDAs. Multiple times in the construction, we will be faced with the task of *moving* or copying the value from some dimension i to the dimension j in the vector. The following fact shows how this can be done using simple matrix multiplication with a specific matrix.

Fact 5.1.5 (Copying elements of a vector). *Let $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{M} \in \mathbb{R}^{D \times D}$ such that $\mathbf{M}_{i,:} = \mathbf{e}_i$, where $\mathbf{M}_{i,:}$ denotes the i^{th} row of \mathbf{M} and \mathbf{e}_i denotes the i^{th} basis vector. Then, it holds that $(\mathbf{M}\mathbf{x})_i = x_i$.*

Also, note that setting the row $\mathbf{M}_{i,:}$ to the zero vector $\mathbf{0}_{\mathbb{R}^D}$ sets the entry x_i to 0, i.e., it erases the entry.

Furthermore, we will use the saturated sigmoid function multiple times to *detect* whether a number of dimensions of a vector are set to one at the same time. Given the recurrent dynamics of the Elman RNN (cf. Eq. (5.26)), we can perform this check as follows.

Fact 5.1.6 (Detecting the activation of multiple values in the hidden state). *Let σ be the saturated sigmoid from Definition 5.1.27, $m \in \{1, \dots, D\}$, $i_1, \dots, i_m, j \in \{1, \dots, D\}$, $\mathbf{x} \in \mathbb{R}^D$, $\mathbf{b} \in \mathbb{R}^D$, and $\mathbf{M} \in \mathbb{R}^{D \times D}$ such that*

$$\mathbf{M}_{j,i} = \begin{cases} 1 & \text{if } i \in \{i_1, \dots, i_m\} \\ 0 & \text{otherwise} \end{cases}$$

and $b_j = -(m-1)$. Then, it holds that $(\sigma(\mathbf{M}\mathbf{x} + \mathbf{b}))_j = 1$ if and only if $x_{i_k} = 1$ for all $k = 1, \dots, m$.²²

Lastly, we will sometimes have to *turn off* certain dimensions of the hidden state if any of the other dimensions are active. Using the dynamics of Elman RNNs and the saturated sigmoid, this can be done as follows.

Fact 5.1.7 (Turning off dimensions in the hidden state). *Let σ be the saturated sigmoid from Definition 5.1.27, $m \in \{1, \dots, D\}$, $i_1, \dots, i_m, j \in \{1, \dots, D\}$, $\mathbf{x} \in \mathbb{R}^D$, $\mathbf{b} \in \mathbb{R}^D$, and $\mathbf{M} \in \mathbb{R}^{D \times D}$ such that*

$$\mathbf{M}_{j,i} = \begin{cases} -1 & \text{if } i \in \{i_1, \dots, i_m\} \\ 0 & \text{otherwise} \end{cases}$$

and $b_j = 1$. Then, it holds that $(\sigma(\mathbf{M}\mathbf{x} + \mathbf{b}))_j = 0$ if and only if $x_{i_k} = 1$ for some $k = 1, \dots, m$.

²²Note that this is simply a restatement of Fact 5.1.4, which we include here for clarity and to make the connection with the construction that follows clearer.

With these intuitions in mind, we now prove Theorem 5.1.11. Due to the relatively elaborate construction, we limit ourselves to pushdown automata with a two-symbol input alphabet $\Sigma = \{a, b\}$ as well as a two-symbol stack alphabet $\Gamma = \{0, 1\}$. Note, however, that this restriction can be done without the loss of generality, meaning that this is enough to prove the Turing completeness of RNNs in general.²³

Proof. We show this by constructing, for a given deterministic pushdown automaton \mathcal{P} recognizing a deterministic context-free language, an Elman RNN simulating the steps performed by \mathcal{P} .

Let $\mathcal{P} = (\Sigma, Q, \Gamma, \delta, (q_i, \gamma_i), (q_\varphi, \gamma_\varphi))$ be such a deterministic pushdown automaton. We now define the parameters of the RNN $\mathcal{R} = (\Sigma, D, \mathbf{U}, \mathbf{V}, \mathbf{E}, \mathbf{b}_h, \mathbf{h}_0)$ such that the updates to \mathcal{R} 's hidden state will correspond to the configuration changes in \mathcal{P} .

The construction is more involved than the one in Minsky's theorem (cf. Lemma 5.1.3). We, therefore, first intuitively describe the semantics of the different *components* of the hidden state of the RNN. Then, we describe the submatrices of the parameters \mathbf{U} , \mathbf{V} , and \mathbf{b} that *control* these components of the vector. The hidden state \mathbf{h} of the RNN will altogether have *five* components.

- **Component 1: Data component:** This component, consisting of three cells, will contain the actual numerical representation of the stack, STACK , as well as two additional “buffer” cells, BUFF_1 and BUFF_2 , which will be used for intermediate copies of the stack values during the computation of the new state.
- **Component 2: Top of stack component:** This component contains three cells, each corresponding to a flag denoting that (a) the stack is empty (STACK_ε), (b) the top element of the stack is a 0 (STACK_0), or (c) the top element of the stack is a 1 (STACK_1).
- **Component 3: Configuration component:** This component encodes the current configuration of the stack (Component 2) together with the current input symbol. Note that, while we assume that the input PDA works with the two-symbol alphabet $\Sigma = \{a, b\}$, the sequence model defined by the RNN requires an EOS symbol to be able to terminate generation (cf. Eq. (2.44)): \mathcal{R} , therefore, defines the conditional probabilities over the set $\bar{\Sigma} = \{a, b, \text{EOS}\}$. With this, there are nine possible configurations (y, γ) for $\gamma \in \{\varepsilon, 0, 1\}$ and $y \in \{a, b, \text{EOS}\}$, meaning that there are nine cells in this configuration, $\text{CONF}_{\gamma, y}$, each corresponding to one of these configurations.
- **Component 4: Computation component:** This component contains four cells in which the computation of the next value of the stack is computed. There are five cells $\text{OP}_{\text{action}, \gamma}$ because *all* possible actions (PUSH 0, PUSH 1, POP 0, POP 1, and NO-OP) are performed simultaneously, and only the correct one is copied into the data component (Component 1) in the end.
- **Component 5: Acceptance component:** This component contains a single cell, ACCEPT , signaling whether the RNN accepts the string \mathbf{y} after reading in the input \mathbf{y} EOS.

Altogether, the hidden state of \mathcal{R} contains $3 + 3 + 9 + 5 + 1 = 21$ dimensions. The *initial hidden state* \mathbf{h}_0 is a vector with a single non-zero component, whose value is 1: the cell STACK_ε since we assume that the stack of the simulated automaton is empty at the beginning of the execution. We now intuitively describe the dynamics that these components define.

The full update step of the network. The RNN will compute the next hidden state corresponding to the new stack configuration by applying the Elman update rule (cf. Eq. (5.26)) four times to complete four discrete sub-steps of the computation. We first define

$$\mathbf{h}_{t+1}^{(0)} \stackrel{\text{def}}{=} \mathbf{h}_t \quad (5.95)$$

and

$$\mathbf{h}_{t+1}^{(n)} \stackrel{\text{def}}{=} \sigma \left(\mathbf{U} \mathbf{h}_{t+1}^{(n-1)} + \mathbf{V} \mathbf{e}(y_t) + \mathbf{b}_h \right) \quad (5.96)$$

for $n = 1, 2, 3, 4$. Then

$$\mathbf{h}_{t+1} \stackrel{\text{def}}{=} \mathbf{h}_{t+1}^{(4)}. \quad (5.97)$$

²³To simulate an arbitrary Turing machine with a machine with the binary alphabet, we simply have to encode each of the finitely-many symbols of the simulated machine using binary encoding.

Intuitively, each of the four stages of computation of the actual next hidden state “detects” some parts of the pattern contributing to the transition in the pushdown automaton. We describe those patterns next intuitively before talking about the submatrices (or subvectors) of the RNN parameters corresponding to the specific parts that update the individual components of the hidden state.

Data component. The cells of the data component form a queue of three components: the `STACK` cell forms the head of the queue, followed by `BUFF1` and `BUFF2`. The values in the cells are updated at each execution of Eq. (5.96) by moving the currently stored values into the next cell in the queue. By doing so, the entry in `BUFF2` gets discarded. The value of `STACK` is copied from the cells of the *computation* component by *summing* them. We will see later that at any point of the computation (when it matters), only one of the computation components will be non-zero, which means that the summation simply corresponds to copying the non-zero computation component. All these operations can be performed by matrix multiplication outlined in Fact 5.1.5.

Encoding the stack sequence. While we outlined a possible encoding of a stack sequence above, the encoding we use in this construction is a bit different. Remember that for a stack sequence $\gamma \in \Gamma^*$ of length N , the *right-most* symbol γ_N denotes the top of the stack. We encode the stack sequence $\gamma \in \Gamma^*$ as follows:

$$\text{rep}(\gamma_1 \dots \gamma_N) \stackrel{\text{def}}{=} \sum_{n=1}^N \text{digit}(\gamma_n) 10^{N-n-1} \quad (5.98)$$

$$\text{where } \text{digit}(\gamma) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } \gamma = 0 \\ 3 & \text{otherwise} \end{cases}.$$

Example 5.1.16. For example, the stack sequence $\gamma = 00110111$ would be represented with $\text{rep}(00110111) = 0.33313311$. Notice the “opposite orientation” of the two strings: the top of the stack in γ is the right-most symbol, while it is the left-most digit in the numerical representation.

Note that the digits 1 and 3 in the definition of $\text{digit}(\cdot)$ are chosen somewhat arbitrarily—the encoding could also have been chosen differently. Similarly, a different (non-decimal) base could have been chosen.

Top of stack component. As mentioned, the `STACKε` cell of this component is one of the two cells set to 1 in the initial state of the RNN. The individual cells of this component then get updated according to the top symbol on the stack encoded in the `STACK` cell by taking into account how the stack is represented by `STACK`. Specifically, the parameters of the RNN are defined such that $\mathbf{h}_{\text{STACK}_\varepsilon} = 1$ if previously $\mathbf{h}_{\text{STACK}} = 0$, $\mathbf{h}_{\text{STACK}_0} = 1$ if previously $\mathbf{h}_{\text{STACK}} = 0.1\dots$, and $\mathbf{h}_{\text{STACK}_1} = 1$ if previously $\mathbf{h}_{\text{STACK}} = 0.3\dots$.

Configuration component. The cells of the configuration component combine the pattern captured by the top of the stack component with the input symbol at the current time step to activate only the appropriate cell `CONFγ,y`. This can be done by incorporating the information from the top of the stack component with the information about the current input symbol from $\mathbf{V}[\llbracket y_t \rrbracket]$. More precisely, the parameters of \mathcal{R} are set such that $\mathbf{h}_{\text{CONF}_{\gamma,y}} = 1$ if at the previous one of the four sub-steps of the computation of the next hidden state, $\mathbf{h}_{\text{STACK}_\gamma} = 1$ and the input symbol is y .

Computation component. The computation component contains the cells in which the results of all the possible actions on the stack are executed. The parameters of the computation component are set such that, given that the previous stack configuration is $\mathbf{h}_{\text{STACK}} = x_1 x_2 \dots x_N$ and the input symbol is y , cells of the computation component are set as

$$\begin{aligned} \mathbf{h}_{\text{OP}_{\text{POP},\gamma}} &= x_2 \dots x_N \\ \mathbf{h}_{\text{OP}_{\text{PUSH},\gamma}} &= \text{digit}(y) x_1 \dots x_N. \end{aligned}$$

Acceptance component. The cell in the acceptance component is activated if and only if the current input symbol is EOS (denoting the end of the string whose recognition should be determined) and the stack is empty, i.e., the STACK_ε cell is activated.

More precisely, the dynamics described here are implemented by the four steps of the hidden state update as follows (where $\mathbf{h}_t^{(1)} = \mathbf{h}_t$).

- In *phase 1*, the configuration of the stack is determined by setting the top of the stack component in $\mathbf{h}_t^{(2)}$.
- In *phase 2*, the configuration of the stack and the input symbol are combined by setting the configuration component in $\mathbf{h}_t^{(3)}$.
- In *phase 3* all possible operations on the stack are performed in the computation component, and, at the same time, the results of all invalid operations (only one operation is valid at each time step due to the deterministic nature of \mathcal{P}) are zeroed-out in $\mathbf{h}_t^{(4)}$. This is done by setting the entries of the recurrence matrix \mathbf{U} such that only the valid action is not zeroed out.
- In *phase 4* the result of the executed operations (only one of which is non-zero) is copied over to the STACK cell in the hidden state in \mathbf{h}_{t+1} .

Having defined the intuition behind the dynamics of the hidden state updates, we now formally define how the parameters of the RNN are set to enable them. Whenever an entry of a matrix or vector is not set explicitly, it is assumed that it is 0 (that is, we only explicitly set the non-zero values). Again, we define them for each component in turn.

The data component. The values of the parameters in the data component are set as follows.

$$U_{\text{BUFF}_1, \text{STACK}} = 1 \quad (5.99)$$

$$U_{\text{BUFF}_2, \text{BUFF}_1} = 1 \quad (5.100)$$

$$U_{\text{STACK}, \text{OP}_{\text{PUSH}, 0}} = U_{\text{STACK}, \text{OP}_{\text{PUSH}, 1}} = U_{\text{STACK}, \text{OP}_{\text{POP}, 0}} = U_{\text{STACK}, \text{OP}_{\text{POP}, 1}} = 1 \quad (5.101)$$

The first two elements correspond to moving the values to the next element in the data component queue, while the entries in the last row correspond to *summing up* the values from the computation component to move them into the stack cell after the computation has been completed. Note that, of course, the elements of the computation component are *always* summed up and written in the STACK cell, no matter what the values there are. However, the division of the computation of the next hidden state into phases ensures that *when it matters*, i.e., after the third phase, there is only a single computation component that is non-zero, and that one is copied into the STACK component in the fourth computation sub-step. All other parameters (in \mathbf{V} and \mathbf{b}_h) are 0.

The top of the stack component. The parameters setting the top of the stack component are set as follows:

$$U_{\text{STACK}_\varepsilon, \text{STACK}} = -10 \quad (5.102)$$

$$U_{\text{STACK}_0, \text{STACK}} = -10 \quad (5.103)$$

$$U_{\text{STACK}_1, \text{STACK}} = 10 \quad (5.104)$$

$$b_{\text{STACK}_\varepsilon} = 1 \quad (5.105)$$

$$b_{\text{STACK}_0} = 3 \quad (5.106)$$

$$b_{\text{STACK}_1} = -2. \quad (5.107)$$

Other parameters (\mathbf{V}) are 0. The reasoning behind these parameters is the following. The cell STACK contains the numeric encoding of the stack content. We distinguish three cases.

- If the stack is *empty*, $\mathbf{h}_{\text{STACK}} = 0$. Therefore, using the parameters above, the value of the cell STACK_1 after the sub-step update will be 0, while the cells STACK_ε and STACK_0 will be 1 due to the positive bias term. This might not be what you would expect—it might seem like, in this case, this step erroneously signals both an empty stack and a stack whose top component is 0. This, however, is corrected for in the configuration component, as we discuss below.
- If the top of the stack is the symbol 0, $\mathbf{h}_{\text{STACK}} = 0.1 \dots$. This means that $10 \cdot \mathbf{h}_{\text{STACK}} \leq 1$ and, therefore, after the update rule application, $\mathbf{h}_{\text{STACK}_1} = 0$. It is easy to see that the setting of the parameters also implies $\mathbf{h}_{\text{STACK}_\varepsilon} = 0$. However, since $-10 \cdot \mathbf{h}_{\text{STACK}} \geq -2$, we have that $\mathbf{h}_{\text{STACK}_0} = 1$.
- Lastly, if the top of the stack is the symbol 1, $\mathbf{h}_{\text{STACK}} = 0.3 \dots$. Therefore, $10 \cdot \mathbf{h}_{\text{STACK}} \geq 3$, meaning that after the update rule application, $\mathbf{h}_{\text{STACK}_1} = 1$. Again, it is easy to see that the setting of the parameters also implies $\mathbf{h}_{\text{STACK}_\varepsilon} = 0$. On the other hand, since $-10 \cdot \mathbf{h}_{\text{STACK}} \leq -3$, it also holds that $\mathbf{h}_{\text{STACK}_0} = 0$.

The configuration component. The configuration component is composed of the most cells of any component:

$$U_{\text{CONF}_{\gamma,y}, \text{STACK}_{\gamma}} = 1 \text{ for } \gamma \in \{\varepsilon, 0, 1\}, y \in \{\text{EOS}, a, b\} \quad (5.108)$$

$$U_{\text{CONF}_{\gamma,0}, \text{STACK}_{\varepsilon}} = -1 \text{ for } y \in \{\text{EOS}, a, b\} \quad (5.109)$$

$$V_{\text{CONF}_{\gamma,y}, m(y)} = 1 \text{ for } \gamma \in \{\varepsilon, 0, 1\}, y \in \{\text{EOS}, a, b\} \quad (5.110)$$

$$b_{\text{CONF}_{\gamma,y}} = -1 \text{ for } \gamma \in \{\varepsilon, 0, 1\}, y \in \{\text{EOS}, a, b\} \quad (5.111)$$

Here, the first, third, and fourth terms together ensure that the cell $\text{CONF}_{\gamma,y}$ is activated if the current input symbol is a ($V_{\text{CONF}_{\gamma,y}, m(y)}$) and the top of the stack is γ ($U_{\text{CONF}_{\gamma,y}, \text{STACK}_{\gamma}}$). $b_{\text{CONF}_{\gamma,y}}$ ensures that both conditions have to be met. The second term, $U_{\text{CONF}_{0,y}, \text{STACK}_{\varepsilon}}$, on the other hand, takes care of an edge case: as shown above, $b_{\text{STACK}_0} = 0$, which means that STACK_0 is, by default, set to 1. The negative weight $U_{\text{CONF}_{0,y}, \text{STACK}_{\varepsilon}} = -1$ ensures that, if the stack is indeed empty, the effect of this default value is “canceled out”, i.e., the configuration cell is not activated by mistake.

The computation component. This is the most complicated component. The computation components are manipulated with the following parameters:

$$U_{\text{PUSH0}, \text{BUFF}_2} = U_{\text{PUSH1}, \text{BUFF}_2} = \frac{1}{10} \quad (5.112)$$

$$U_{\text{POP0}, \text{BUFF}_2} = U_{\text{POP1}, \text{BUFF}_2} = 10 \quad (5.113)$$

$$U_{\text{NO-OP}, \text{BUFF}_2} = 1 \quad (5.114)$$

$$U_{A, \text{CONF}_{\gamma,y}} = -10 \text{ for } A \in \{\text{OP}_{\text{PUSH},0}, \text{OP}_{\text{PUSH},1}, \text{OP}_{\text{POP},0}, \text{OP}_{\text{POP},1}, \text{OP}_{\text{NO-OP}}\} \\ \gamma \in \{0, 1\}, y \in \{a, b\} \quad (5.115)$$

$$U_{\text{OP}_{A,\gamma'}, \text{CONF}_{\gamma,y}} = 0 \text{ for } q \xrightarrow{y, \gamma \rightarrow \gamma'} q \in \delta, \\ A \in \{\text{OP}_{\text{PUSH},0}, \text{OP}_{\text{PUSH},1}, \text{OP}_{\text{POP},0}, \text{OP}_{\text{POP},1}, \text{OP}_{\text{NO-OP}}\} \quad (5.116)$$

$$U_{\text{OP}_{\text{NO-OP}}, \text{CONF}_{\gamma, \text{EOS}}} = 0 \text{ for } \gamma \in \{\varepsilon, 0, 1\} \quad (5.117)$$

$$b_{\text{OP}_{\text{PUSH},0}} = \frac{1}{10} \quad (5.118)$$

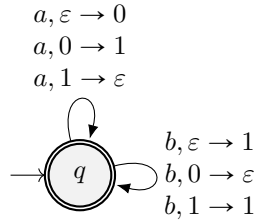
$$b_{\text{OP}_{\text{PUSH},1}} = \frac{3}{10} \quad (5.119)$$

$$b_{\text{OP}_{\text{POP},0}} = -1 \quad (5.120)$$

$$b_{\text{OP}_{\text{POP},1}} = -3 \quad (5.121)$$

$$b_{\text{OP}_{\text{NO-OP}}} = 0 \quad (5.122)$$

The first three parameters above concern *copying* the value of the stack encoded by the previous hidden state into the computation component and preparing it for modification. They work together with the corresponding entries in the bias vector \mathbf{b}_h . For example, a value can be *pushed* onto the stack by dividing the

Figure 5.16: The single-stack pushdown automaton \mathcal{P} .

value of the stack encoding by 10 and adding either 0.1 or 0.3, depending on whether 0 or 1 is being pushed. This is encoded by the first setting above and $b_{\text{OP}_{\text{PUSH},0}}$ and $b_{\text{OP}_{\text{PUSH},1}}$. Similarly, a value can be popped from the stack by multiplying the stack encoding with 10 and then subtracting the appropriate value according to the bias entry. The NO-OP action is implemented simply by copying the values of the stack into its cell. The remaining three parameter settings above ensure that, after executing all possible stack actions, *only the appropriate computation* is kept and all others are zeroed out. The fourth row above ensures that “by default”, all computation cells are reset to 0 after every update. However, the next row “removes” the negative weights (sets them to 0) for the changes in the configuration which correspond to the *valid transitions*, or valid actions, in the pushdown automaton. That is, setting those values of the matrix \mathbf{U} to 0 disables “erasing” the entry $\text{OP}_{A,\gamma'}$ in the hidden state by the configuration $\text{CONF}_{\gamma,y}$ *if* the transition from the configuration with the top of the stack γ to γ' with the action A upon reading y is encoded by the original automaton. The last remaining row simply ensures that reading in the EOS symbol results in the NO-OP action being executed (EOS actions are not encoded by the original pushdown automaton).

The acceptance component. Lastly, the acceptance component is controlled by the following parameters:

$$U_{\text{ACCEPT},A} = -10 \text{ for } A \in \{\text{OP}_{\text{PUSH},0}, \text{OP}_{\text{PUSH},1}, \text{OP}_{\text{POP},0}, \text{OP}_{\text{POP},1}, \text{OP}_{\text{NO-OP}}\} \quad (5.123)$$

$$U_{\text{ACCEPT},\text{CONF}_{\gamma,y}} = -10 \text{ for } \gamma \in \{\varepsilon, 0, 1\}, y \in \{\text{EOS}, a, b\} \quad (5.124)$$

$$b_{\text{ACCEPT}} = 1 \quad (5.125)$$

The entry b_{ACCEPT} ensures that, by default, the value of ACCEPT is set to 1. However, the other parameters ensure that, as soon as any part of the configuration is not compatible with the acceptance state (the read symbol is not EOS or the stack is not empty), the acceptance bit is turned off.

A full proof of the theorem would now require us to show formally that the update rule Eq. (5.97) results in the correct transitions in the PDA. We, however, leave the proof with the intuitive reasoning behind the setting of the parameters and leave this as an exercise for the reader. The proof is also demonstrated in the python implementation of the constructions here: <https://github.com/rycolab/rnn-turing-completeness>. ■

The construction described in the proof of Theorem 5.1.11 is demonstrated in the following example.

Example 5.1.17. Let \mathcal{P} be a single-stack PDA presented in Fig. 5.16. We now simulate the recognition of the string $y = ab$, which is accepted by \mathcal{P} . The initial state \mathbf{h}_0 has a single non-zero cell, STACK_ε . The four phases of the processing of the first input symbol a are shown in Tab. 5.1. The four phases of the processing of the second input symbol b are shown in Tab. 5.2.

Theorem 5.1.11 shows that Elman RNNs are theoretically at least as expressive as deterministic CFGs. We now return to the main result of this subsection: the Turing completeness of RNNs. Luckily, Theorem 5.1.11 gets us most of the way there! Recall that by Theorem 4.2.6, *two-stack* PDA are Turing complete. We make use of this fact by generalizing the construction in the proof of Theorem 5.1.11 to the two-stack case. This will prove that RNNs can in fact simulate any Turing machine, and are, therefore, Turing complete.

Lemma 5.1.18. Let $\mathcal{P} = (Q, \Sigma, \Gamma, \delta, (q_i, \gamma_i, \sigma_i), (q_\varphi, \gamma_\varphi, \sigma_\varphi))$ be a two-stack pushdown automaton. Then, there exists an Elman RNN \mathcal{R} simulating \mathcal{P} , i.e., $L(\mathcal{R}) = L(\mathcal{P})$.

	Initial state	Phase 1	Phase 2	Phase 3	Phase 4
STACK	0	0	$\frac{2}{5}$	$\frac{2}{5}$	$\frac{1}{10}$
BUFF ₁	0	0	0	$\frac{3}{5}$	$\frac{2}{5}$
BUFF ₂	0	0	0	0	$\frac{1}{10}$
STACK _ε	1	1	1	0	0
STACK ₀	0	1	1	0	0
STACK ₁	0	0	0	1	1
CONF _{EOS,a}	0	0	1	1	0
CONF _{EOS,b}	0	0	0	0	0
CONF _{0,a}	0	0	0	0	0
CONF _{0,b}	0	0	0	0	0
CONF _{1,a}	0	0	0	0	1
CONF _{1,b}	0	0	0	0	0
CONF _{ε,EOS}	0	0	0	0	0
CONF _{0,EOS}	0	0	0	0	0
CONF _{1,EOS}	0	0	0	0	0
OP _{PUSH,0}	0	$\frac{1}{10}$	$\frac{1}{10}$	$\frac{1}{10}$	$\frac{1}{10}$
OP _{PUSH,1}	0	$\frac{3}{10}$	$\frac{3}{10}$	0	0
OP _{POP,0}	0	0	0	0	0
OP _{POP,1}	0	0	0	0	0
OP _{NO-OP}	0	0	0	0	0
ACCEPT	0	1	0	0	0

Table 5.1: The simulation of the processing of the first symbol a by the RNN simulating the PDA in Fig. 5.16. After the fourth phase, the stack cell contains the encoding of the stack as 0.1.

	Initial state	Phase 1	Phase 2	Phase 3	Phase 4
STACK	1/10	1/10	1	2/5	0
BUFF ₁	2/5	1/10	1/10	1	2/5
BUFF ₂	2/5	2/5	1/10	1/10	1
STACK _ε	0	0	0	0	0
STACK ₀	0	1	1	0	0
STACK ₁	1	0	0	1	1
CONF _{EOS,a}	0	0	0	0	0
CONF _{EOS,b}	0	0	0	0	0
CONF _{0,a}	0	0	0	0	0
CONF _{0,b}	0	0	1	1	0
CONF _{1,a}	1	0	0	0	0
CONF _{1,b}	0	1	0	0	1
CONF _{ε,EOS}	0	0	0	0	0
CONF _{0,EOS}	0	0	0	0	0
CONF _{1,EOS}	0	0	0	0	0
OP _{PUSH,0}	1/10	0	0	0	0
OP _{PUSH,1}	0	0	0	0	0
OP _{POP,0}	0	0	0	0	0
OP _{POP,1}	0	1	0	0	0
OP _{NO-OP}	0	0	2/5	0	0
ACCEPT	0	0	0	0	0

Table 5.2: The simulation of the processing of the second symbol b by the RNN simulating the PDA in Fig. 5.16. After the fourth phase, the stack cell contains the encoding of the empty stack.

Proof. Again, given a two-stack PDA $\mathcal{P} = (Q, \Sigma, \Gamma, \delta, (q_\iota, \gamma_\iota, \sigma_\iota), (q_\varphi, \gamma_\varphi, \sigma_\varphi))$ with $\Sigma = \{a, b\}$, $\Gamma_1 = \{0, 1\}$, and $\Gamma_2 = \{0, 1\}$, we construct an Elman RNN $\mathcal{R} = (\Sigma, D, \mathbf{U}, \mathbf{V}, \mathbf{E}, \mathbf{b}_h, \mathbf{h}_0)$ which recognizes the same language as \mathcal{P} .

The hidden state of \mathcal{R} will contain the same components as in the proof of Theorem 5.1.11. Moreover, their dynamics will be exactly the same—they will simply be *larger* to account for more possible configurations of the two stacks together. For example, the top of the stack component will now consist of cells $\text{STACK}_{\gamma_1\gamma_2}$ for $\gamma_1 \in \Gamma_1$ and $\gamma_2 \in \Gamma_2$ flagging that the top symbol on stack 1 is γ_1 and the top symbol of stack 2 is γ_2 . Furthermore, the configuration component will contain cells of the form $\text{OP}_{\text{action}, \gamma_1\gamma_2}$ with an analogous interpretation. Lastly, all the computation and data component cells would be duplicated (with one sub-component for each of the stacks), whereas the acceptance component stays the same. We now again describe these components and their dynamics intuitively, whenever there is any difference to the single-stack version.

Data component. Instead of having a *single* queue of data cells, \mathcal{R} now has *two* queues, one for each of the stacks. The first queue will be formed by STACK_1 , BUFF_{11} , and STACK_{21} , and the second one by STACK_2 , BUFF_{12} , and STACK_{22} . Each of the queues acts exactly the same as in the single-stack version, and they act independently based on the computations done in the computation cells of the respective stacks. Each of the stacks is also encoded as a numeric sequence in the same way as in the single-stack version.

Top of stack component. Again, \mathcal{R} starts in an initial state in which the cell $\text{STACK}_{\varepsilon\varepsilon}$ is 1 and all others are 0. The individual cells of this component then get updated according to the top symbols on the stacks encoded in the STACK_1 and STACK_2 cells.

Configuration component. The cells of the configuration component combine the pattern captured by the top of *both* stack components with the input symbol at the current time step to activate only the appropriate cell $\text{CONF}_{\gamma_1\gamma_2, y}$.

Computation component. Again, the computation component contains the cells in which the results of all the possible actions on both stacks are executed. They execute the actions on both stacks independently.

Acceptance component. The acceptance component functions identically to the single-stack case.

Using these components, the RNN then transitions between the phases exactly like in the single-stack case. We leave the specifications of the matrix parameter values to the reader. They again follow those presented in the single-stack case but treat the transitions and configurations of both stacks. ■

The Computational Power of RNN Variants

Most of the current section was devoted to understanding the computational power of simple Elman RNN language models due to their simplicity which allows an easy connection to formal models of computation. However, as mentioned in §5.1.6, gated RNN variants such as LSTMs and GRUs have become the standard in modern natural language processing tasks, showing more better and more reliable performance on a variety of tasks. Interestingly, besides their better resilience to the vanishing gradient problem, LSTM-based language models are also provably more *expressive* than simple RNN language models. On the other hand, the simpler GRU-based language models are in some ways only as powerful as simple Elman RNN language models. To give an intuition behind this, this subsection provides a short overview the results considering the computational power of RNN variants.

Weiss et al. (2018) compare the practical computational power of different RNN variants under the constraint of bounded computation time and limited precision. Interestingly they empirically find that LSTMs can learn to recognize languages that require some form of *counting*, like $\{a^n b^n \mid n \in \mathbb{N}\}$ or $\{a^n b^n c^n \mid n \in \mathbb{N}\}$, while GRUs struggle to do so. This invites the comparison of LSTMs with **counter machines** (Hopcroft et al., 2006; Fischer et al., 1968), a class of formal computational models with the ability to count. Simply put, counter machines are finite-state automata with an additional (unbounded) counter cell, which they

can manipulate by incrementing and decrementing the value stored in it.²⁴ Counter machines present an interesting addition to the traditional hierarchy of computational models, since they in some ways *cut across* it, being able to recognize some, but not all, context-free languages, while also being able to recognize some *context-sensitive* languages. Among others, they can for example recognize the context-free language $\{a^n b^n \mid n \in \mathbb{N}\}$ and the context-sensitive language $\{a^n b^n c^n \mid n \in \mathbb{N}\}$, while not being able to recognize the Dyck context-free languages $D(k)$ with k different parenthesis types (intuitively, this is because recognizing Dyck languages requires keeping track of the order in which the parantheses appeared, not only their counts).

Counter machines can recognize languages like $\{a^n b^n c^n \mid n \in \mathbb{N}\}$, by counting the number of *as* and making sure that it matches the number of *bs*. Further, analyzing the activation of the memory cell and of the hidden state, they find that LSTMs that one can recognize the use of a counting mechanism implemented by the network by using one or more dimensions of the memory cell as counters. This result is particularly interesting as GRUs are often considered an equivalent variant to the LSTM one, with the same computational power, but smaller computation overhead. However, GRUs seem to lack this counting behavior, which can be backed by theoretical analysis.

Merrill et al. (2020) take a look at the computational power of the different RRN variant from another perspective. They consider space complexity and whether the networks are rationally recurrent, i.e. whether their hidden state update function can be expressed in terms of finite state machine computations. Making the assumption of saturated networks²⁵, they find that while GRUs and Elman networks are rationally recurrent and therefore at most regular, LSTMs are not, meaning that their hidden state update function cannot be expressed by means of finite-state machines.

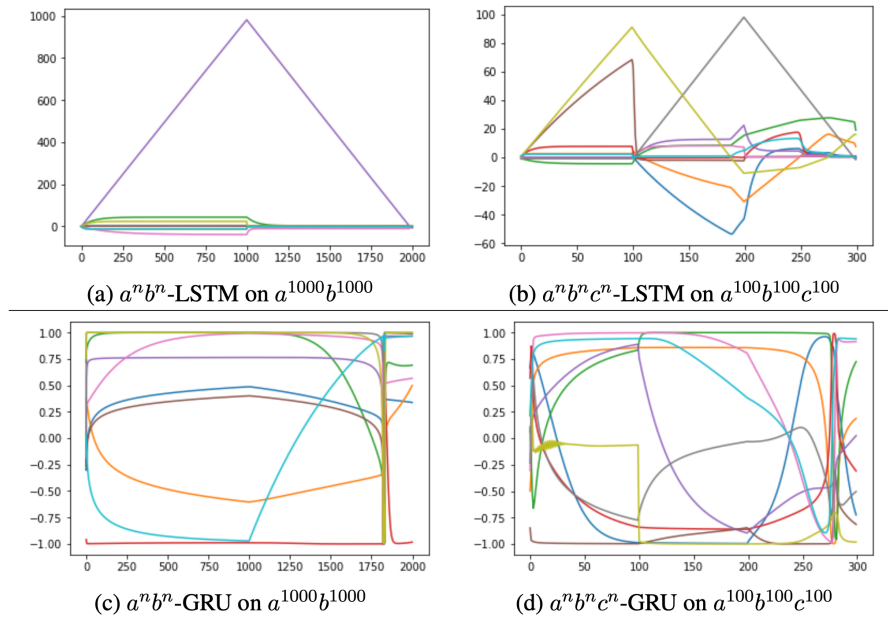


Figure 5.17: Picture from Weiss et al. (2018): Plot of the activation value of the memory cell (LSTM) and of the hidden state (GRU), versus the indices of the input. The Networks have been trained to recognize either the languages $\{a^n b^n \mid n \in \mathbb{N}\}$ or $\{a^n b^n c^n \mid n \in \mathbb{N}\}$. As you can notice, in both cases the LSTM has learned to use one or more dimension of the memory cell as a *counter*, which allows it to count how many *a* and *b* have been consumed so far. Conversely, the GRU hasn't developed this counter mechanism, and in fact empirical evidence shows that it struggles to recognize the described languages.

²⁴We only consider counter machines with a *single* counter cell.

²⁵Informally, a saturated network is a neural network in which the norms of the parameters are taken to ∞ . This has the effect of pushing all the squashing functions of the network to one of their extreme values, in the case of the sigmoid $\{0, 1\}$ and $\{-1, 1\}$ in the case of the hyperbolic tangent

Consequences of the Turing completeness of recurrent neural networks

The section above outlines Siegelmann and Sontag’s (1992) construction of encoding a Turing machine in an RNN. While Turing completeness means that RNNs are in many ways computational very powerful (as powerful as they can be), it also brings with it many computational challenges faced when working with Turing machines. Computability theory, for example, defines many problems related to Turing machines and their properties which are *not computable*, or *undecidable*, meaning that no algorithm (or, equivalently, a Turing machine) can solve them.²⁶ The most classical and fundamental of such problems is the halting problem (Turing, 1937).

Definition 5.1.28. *Let \mathcal{M} be a Turing machine over the input alphabet Σ and $\mathbf{y} \in \Sigma^*$. The **halting problem** is the problem of deciding whether \mathcal{M} (halts and) accepts \mathbf{y} .²⁷*

The halting problem is the foundation of the results presented by Chen et al. (2018), which building on Siegelmann and Sontag (1992) considers its implications on the theory of RNN language models. For example, they show that determining many practically useful properties of general rationally weighted RNNs is *undecidable*. Such properties include the tightness of a general RNN,²⁸ the equivalence of two RNN language models, the minimal size (in the size of the hidden state) of an RNN defining the same language model as some given RNN, and the highest probability string in an RNN, i.e., $\arg\max_{\mathbf{y} \in \Sigma^*} p_{\text{LN}}(\mathbf{y})$. By simulating a Turing machine, we can encode the halting problem as solving any of those tasks.²⁹ This means that if we could solve these problems for RNNs, we could also solve the halting problem, which is provably impossible, meaning that these problems are not computable either. We briefly outline the individual findings of Chen et al. (2018) below, sketching the rough idea behind the proofs.³⁰

Theorem 5.1.12 (Tightness). *Determining tightness of an RNN language model is undecidable.*

Proof. Note first that not all RNNs are tight language models. As mentioned, this is not a contradiction to the results from §5.1.3, which considered softmax-projected RNN language models. Indeed, the RNN we constructed in §5.1.7 is not such an RNN. That construction shows that, given an arbitrary Turing machine \mathcal{M} and input \mathbf{y} , we can construct an RNN that simulates \mathcal{M} running on \mathbf{y} . Using this construction, we can reduce the problem of deciding the tightness of a general (rational weighted Elman) RNN to the halting problem, i.e., the problem of answering “Does \mathcal{M} halt on \mathbf{y} ?”. We do that by constructing an RNN that simulates the given Turing machine \mathcal{M} on the input \mathbf{y} , ending generation if \mathcal{M} halts, and, at the same time, produces strings according to a distribution that is tight for finite length strings, on the condition that no infinite length strings can be produced. Now if and only if \mathcal{M} halts on \mathbf{y} , the language model is tight. Therefore, deciding is at least as hard as solving the halting problem, and since that is undecidable, so is tightness. ■

Theorem 5.1.13 (Highest weighted string). *Finding the highest weighted string of an RNN LM is undecidable.*

Proof. Once more, we reduce the halting problem to the given task on the RNN. The idea behind this is to again simulate an arbitrary Turing machine by constructing an RNN LM which is not tight unless the simulated Turing machine halts. One can create such an RNN with a one-symbol output alphabet such that its highest-weighted output is an infinite string. Now, if we again enforce that the RNN ends producing outputs once the simulated Turing machine halts, it can be shown that there exists a language in which each string has a probability of less than 0.12 if and only if the Turing machine does not halt. On the other hand, if the Turing machine does halt after T steps, producing a string which has length $3T - 5$ has a probability of ≥ 0.25 . Therefore, the weight of the highest probability string depends on the question of whether the simulated Turing machine halts, which is undecidable. ■

²⁶The notion of solving a problem computationally is quite nuanced and we only provide very broad intuitions for now. Readers who want to delve deeper into the topic of computability and with it the implications of Turing completeness are encouraged to look at some classical textbooks on the material, for example Sipser (2013, Part Two).

²⁷A formal definition of the halting problem would require us to define the formal language L of Turing machine and input tuples $(\mathcal{M}, \mathbf{y})$ and asking if there exists a Turing machine \mathcal{M}' that accepts L . However, to keep the discussion brief, we define the problem more intuitively.

²⁸Note that the results from §5.1.3 consider only RNN LM with the softmax projection function.

²⁹In computability theory, this is known as *reducing* the halting problem to the given one

³⁰The interested reader is encouraged to see the original paper for the details.

Theorem 5.1.14 (Equivalence). *Equivalence between two RNN LMs in the sense of defining the same language model is undecidable.*

Proof. The proof of this claim is again a reduction from the halting problem. We construct an RNN which simulates a given arbitrary Turing machine until it halts and has the same outputs as some other RNN. As soon as the Turing machine halts, the outputs differ, so the RNNs are equivalent if and only if the Turing machine does not halt. Hence, equivalence is undecidable. ■

Theorem 5.1.15 (Minimization). *Finding the RNN with the minimum number of hidden layer neurons defining the same language model as a given RNN LM is undecidable.*

Proof. We can reduce the halting problem to the following problem: For some RNN LM and an integer D , return yes if there is another RNN LM with $\leq D$ hidden units that generates the same weighted language. Assume that there is a Turing machine \mathcal{M} that can decide this problem. Now, for another Turing machine \mathcal{M}' and input \mathbf{y} , construct a one symbol RNN LM, \mathcal{R} , that simulates \mathcal{M}' running on \mathbf{y} and stops generating if \mathcal{M}' halts. We assume without loss of generality that \mathcal{M}' runs for more than one computation step. Now we run \mathcal{M} on the input $(\mathcal{R}, 0)$, which checks whether there is another RNN LM generating the same weighted language as \mathcal{R} and has no hidden units. Having no hidden units means that the output probabilities of each symbol would have to be constant for each time step. If \mathcal{M} decides minimization returns true, that means the output probabilities of \mathcal{R} cannot change over time, which means that \mathcal{M}' has to run forever. Conversely, if \mathcal{M} returns false, the output probabilities change when \mathcal{M}' halts. Therefore, \mathcal{M} deciding the minimal hidden states problem on $(\mathcal{R}, 0)$ is equivalent to it deciding the Halting problem for $(\mathcal{M}', \mathbf{y})$. ■

This concludes our investigation of the formal properties of recurrent neural language models. The sequential nature of the architecture and the relatively simple transition functions in the vanilla RNN architectures made the link to automata from formal language theory relatively straightforward, which allowed relatively strong theoretical insights. However, it was exactly this sequential nature and the issues associated with it of RNNs that eventually led to them being overtaken by another neural architecture, which is now at the core of most if not all, modern state-of-the-art language models: the transformer.³¹ We introduce them and discuss their theoretical underpinnings in the next section.

³¹We will not discuss the issues with the training speed and parallelization of RNNs in detail. Some of these issues will be highlighted in the latter parts of the course.

5.2 Transformer-based Language Models

In the previous section (§5.1.2), we introduced and studied RNN language models as of language models capable of storing an arbitrarily long context in its encoding $\text{enc}(\mathbf{y}_{<t})$ by updating their hidden state \mathbf{h}_t an arbitrary number of times. As we saw in their theoretical analysis, this mechanism gives them, under some assumptions, a lot of expressive power. Nevertheless, as we also discussed, RNN LMs come with their distinct set of drawbacks. Some of them, e.g., the exploding and vanishing gradient problems, can be amended using specific mechanisms resulting in more complex recurrent neural networks, such as LSTMs and GRUs (cf. §5.1.6). As discussed in §5.1.6, a more fundamental issue that cannot be avoided is the difficulty of parallel training, which is particularly noticeable on the vast internet-scale corpora used nowadays to train language models. Can we do anything about that? As discussed, the inherently sequential nature of RNNs suggests strict limits on this front. Motivated by this limitation, in this section, we present a newer architecture that took over the field of language modeling (and NLP in general)—transformers (Vaswani et al., 2017). It was originally introduced for machine translation, but it can easily be applied to language modeling and has led to the success of models such as GPT-n.

The structure of this section will be a bit different from that of the other sections in the notes. We will first give a formal definition of a transformer model in §5.2.2 and, based on this definition, derive a number of results analogous to those for RNN LMs from §5.1.7. However, due to the current practical relevance of transformers in language modeling, we then devote a significant portion of the section to more practical aspects of transformer models and introduce a number of modifications used in modern language modeling systems.

5.2.1 Informal Motivation of the Transformer Architecture

Before we introduce transformers, let us consider another practical drawback of RNN LMs, which will give us more clues on how to improve them and motivate the architectural decisions behind transformers. Luckily, the simplest patch-up of this issue will also lend itself naturally to parallelization, as we will see shortly. The main characteristic of RNNs is the use of a single hidden state \mathbf{h}_t to represent an arbitrary prefix of any string $\mathbf{y}_{<t}$ up to the current time step t . While this allows RNNs to model strings of any length, it also means that arbitrarily long strings must be *compressed* into this hidden vector of fixed size. Intuitively, this becomes increasingly difficult as the length of the context grows: As the amount of information to be compressed into the hidden state increases with the prefix length, the hidden state may struggle to model the entirety of the preceding context. How can we amend that? The simplest naïve way to go about this is to retain the contextual encodings of *all* prefixes of the string. In this case, it is actually more natural to talk about contextual encodings not of full prefixes, but simply of the individual symbols in the string.³² Here, contextual means that the symbol encodings are augmented with the information from the rest of the string (in most cases, about the preceding context, as with the hidden states of an RNN). With this, we avoid the need to summarize the entire context into a single state. Note that our infinite-precision RNNs from the previous section implicitly did that as well—for example, by storing the information in the “stack” neurons, they could, in principle, store the entire history of the string. However, storing all the encodings explicitly makes their utilization more direct and thus easier. This of course leaves the model with the issue of remembering increasingly large amounts of information as the length of the context increases, but we will, for the moment, assume that we can always store enough information to process any string in this way.

Having decided to keep around the encodings of all symbols in the string, let us think about parallelizing the process of encoding a string, i.e., computing $\text{enc}(\mathbf{y})$. Remember the very general way in which RNNs build a representation of the string $\mathbf{y}_{<t}$ by incrementally modifying \mathbf{h}_t , which is illustrated in Fig. 5.1a—this incrementality brings with it all the challenges of impossible parallelization. The workaround for the issues with the sequential processing of RNNs is to process the context for each y_t *independently*, without relying on the encodings of the previous symbols, thus avoiding the sequential bottleneck. Nevertheless, we still want the contextual encoding of y_t to contain information about the rest of the string, i.e., the preceding context. How can we achieve that without relying on recurrence? Again, we grab onto the simplest solution: to compute the symbol encodings for each symbol y_t “from the ground up” based only on the static symbol encodings

³²For example, looking back to RNNs, we could consider \mathbf{h}_t to simply be an encoding of the symbol y_t augmented with the information from $\mathbf{y}_{<t}$.

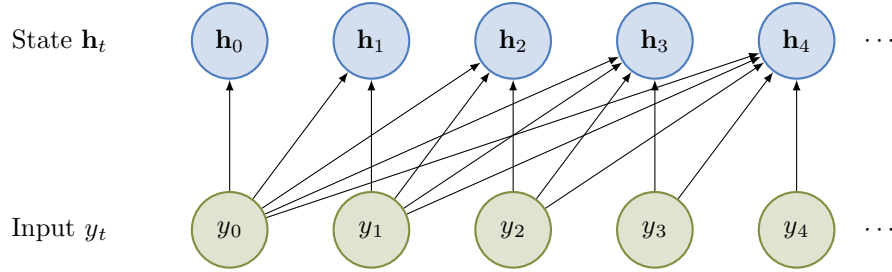


Figure 5.18: An abstract depiction of how a transformer language model produces the contextual embeddings of all symbols in a string. The hidden state \mathbf{h}_t can “attend to” (the precise meaning of this term will be introduced soon in §5.2.2) all preceding symbols $\mathbf{y}_{<t}$ and the current symbol \mathbf{y}_t .

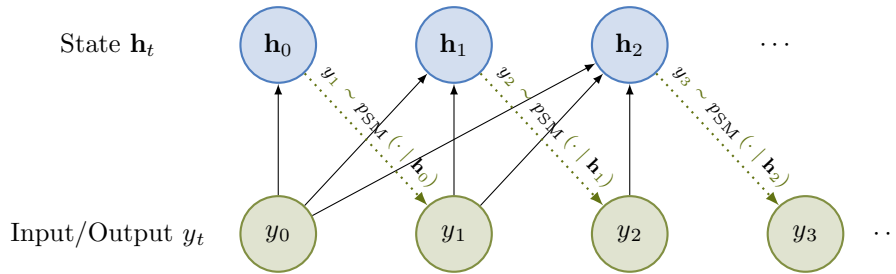


Figure 5.19: An abstract depiction of how a transformer language model *generates* a string one symbol at a time. The hidden state \mathbf{h}_t can attend to all previously generated symbols $\mathbf{y}_{<t}$ to sample the next symbol \mathbf{y}_t . The dotted lines denote the sampling steps.

$\mathbf{e}'(y_t)$, which do not require any recurrence. This is abstractly illustrated in Fig. 5.18, whereas Fig. 5.19 shows how this translates into the generative framework from Definition 3.1.11, where individual symbols y_t are both sequentially generated based on the encodings of the preceding context $\text{enc}(\mathbf{y}_{<t})$ as well as used to build the representation of the context in the next time step $\text{enc}(\mathbf{y}_t)$. Notice that instead of containing arcs denoting dependencies between the symbol encodings (“hidden states”) \mathbf{h}_t , Fig. 5.18 and Fig. 5.19 contain arcs connecting each \mathbf{h}_t to all symbols y_j for all $j \leq t$. Compare this to Fig. 5.1a, where the arcs between \mathbf{h}_{t-1} and \mathbf{h}_t induce the temporal dependence, and carry the information about the symbols y_j for all $j \leq t$.

Clearly, this avoids the issues faced by RNNs due to their sequentiality. However, it also introduces more work required to compute the individual contextual encodings from the static symbol representations. The operations performed to do this by transformers, which are together known as the **attention mechanism**, are introduced in the next subsection. They represent possibly the most important component of the entire structure of the transformer—by “attending” to relevant preceding symbols when computing the symbol encodings (i.e., using them to compute $\text{enc}(\mathbf{y}_{<t})$), the transformer can model long-range dependencies very effectively, and use them for appropriately modeling the distribution over the next word.

To recap, in this section, we informally motivated the new architecture, transformers, with the goals of 1. remembering the contextual encodings of all symbols explicitly and 2. parallelizing the computation of the contextual symbol encodings. The next subsection formally introduces the architecture, before we dive into their theoretical properties.

5.2.2 A Formal Definition of Transformers

Having informally introduced the main two ideas behind the transformer architecture in the previous subsection, we now provide a formal definition of a transformer model, which we will then augment with more practical considerations in .

Definition 5.2.1 (Transformer network). *A **transformer network** \mathcal{T} is a tuple $(\Sigma, D, \text{enc}_{\mathcal{T}})$ where*

- Σ is the alphabet of input symbols,

- D is the dimension of \mathcal{T} , and
- $\text{enc}_{\mathcal{T}}$ is the transformer encoding function (cf. Definition 3.1.7), which we define in more detail below (Definition 5.2.2).

From afar, the definition of a transformer network is therefore relatively simple; it is stated to make the transformer models fit well into our representation-based locally normalized language modeling framework (cf. Definition 3.1.11). The complexity of the models of course comes from the definition of the transformer encoding function $\text{enc}_{\mathcal{T}}$, to which we devote the rest of the section.

Continuing in the framework of representation-based locally normalized language models, the hidden states of the transformer play an analogous role to those of RNNs, with the only difference being how they are computed.

Definition 5.2.2 (Transformer hidden state). *Let $\mathcal{T} = (\Sigma, D, \text{enc}_{\mathcal{T}})$ be a transformer network. The hidden state $\mathbf{h}_t \in \mathbb{R}^D$ describes the state of \mathcal{T} after reading $\mathbf{y}_{\leq t}$. It is defined with respect to the transformer encoding function $\text{enc}_{\mathcal{T}}$ as follows:*

$$\mathbf{h}_t \stackrel{\text{def}}{=} \text{enc}_{\mathcal{T}}(\mathbf{y}_{\leq t}) \quad (5.126)$$

Crucially, as we will see shortly, the hidden state \mathbf{h}_t of the transformer does not have any dependence on the preceding hidden states themselves (although, as we will see, it is partially a function of the same *inputs*).

As hinted above, with this, we can easily fit transformers into the representation-based locally normalized language modeling framework and define a sequence model based on the model.

Definition 5.2.3 (Transformer sequence model). *Let \mathcal{T} be a transformer network and $\mathbf{E} \in \mathbb{R}^{|\Sigma| \times D}$ a symbol representation matrix. A D -dimensional **transformer sequence model** over the alphabet Σ is a tuple $(\Sigma, D, \text{enc}_{\mathcal{T}}, \mathbf{E})$ defining the sequence model of the form*

$$p_{SM}(\bar{y}_t \mid \mathbf{y}_{<t}) \stackrel{\text{def}}{=} \text{softmax}(\mathbf{E} \mathbf{h}_{t-1})_{\bar{y}_t} = \text{softmax}(\mathbf{E} \text{enc}_{\mathcal{T}}(\mathbf{y}_{<t}))_{\bar{y}_t} \quad (5.127)$$

Now that we have unified the transformer \mathcal{T} to the theoretical framework introduced so far in the course, we can jump in and look at the internal structure of the transformer encoder function, which is where the novelty of the transformer architecture comes from.

The Attention Mechanism

As we mentioned in the informal motivation, to avoid over-compressing information about sentences into a single vector, a transformer model retains the encodings (captured in the hidden states \mathbf{h}_t) of *all* possible prefixes of the string, which we can equivalently simply regard as encodings of individual symbols augmented with the information from the preceding string (see Fig. 5.18).³³ However, rather than computing the encodings sequentially like an RNN, the encodings of the individual symbols are computed with the so-called attention mechanism.

Definition 5.2.4 (Attention). *Let $f : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ be a **scoring function** and $\mathbf{f}_{\Delta^{D-1}}$ a projection function. Furthermore, let $\mathbf{q} \in \mathbb{R}^D$, $\mathbf{K}_t = (\mathbf{k}_1^\top, \dots, \mathbf{k}_t^\top) \in \mathbb{R}^{t \times D}$ and $\mathbf{V}_t = (\mathbf{v}_1^\top, \dots, \mathbf{v}_t^\top) \in \mathbb{R}^{t \times D}$.*

Attention over $\mathbf{K}_t, \mathbf{V}_t$, also denoted by $\text{Att}(\mathbf{q}, \mathbf{K}_t, \mathbf{V}_t) : \mathbb{R}^D \times \mathbb{R}^{t \times D} \times \mathbb{R}^{t \times D} \rightarrow \mathbb{R}^D$ is a function computing the vector \mathbf{a} in the following two-step process:

$$\mathbf{s}_t = (s_1, \dots, s_t) \stackrel{\text{def}}{=} \mathbf{f}_{\Delta^{D-1}}(f(\mathbf{q}, \mathbf{k}_1), f(\mathbf{q}, \mathbf{k}_2), \dots, f(\mathbf{q}, \mathbf{k}_t)) \quad (5.128)$$

$$\mathbf{a}_t = \text{Att}(\mathbf{q}, \mathbf{K}_t, \mathbf{V}_t) \stackrel{\text{def}}{=} s_1 \mathbf{v}_1 + s_2 \mathbf{v}_2 + \dots + s_t \mathbf{v}_t \quad (5.129)$$

\mathbf{q} , \mathbf{K} , and \mathbf{V} are commonly referred to as the **query**, **keys**, and **values** of the attention mechanism, respectively. We talk about the parameters \mathbf{q} , \mathbf{K} , and \mathbf{V} completely abstractly for now. However, to help you connect this to the representation-based language modeling framework, note that \mathbf{q} will later correspond to a query representing an individual symbol y_t , whereas \mathbf{K} and \mathbf{V} will contain the information from $\mathbf{y}_{<t}$ used to compute \mathbf{h}_t .

³³From now on, we will talk about **contextual symbol encodings**, which simply refers to the hidden states corresponding to individual symbols.

What the attention function computes. The scoring function f is, abstractly, simply a parameter of the model which we can choose freely. Intuitively, it should express the relevance of a particular key \mathbf{k} to the query \mathbf{q} —the more the key is relevant to the query, the more “attention” the model will put to the *value* associated to that key. The projection function $\mathbf{f}_{\Delta^{D-1}}$ then transforms the computed scores ensuring that the transformed scores sum to 1.³⁴ The vector of transformed scores \mathbf{s} (Eq. (5.128)) is then used to compute the result of the attention function—the vector \mathbf{a} . \mathbf{a} is a convex combination of the *values* \mathbf{v} passed to the attention function. Abstractly, therefore, the *keys* contain the information used for “indexing” the *values* with the specific *query*.

The scoring function. As mentioned, the scoring function is supposed to measure the “relevance” of a particular value for a query \mathbf{q} through the values’ key. The most common choice for f is the dot product between query and key, which is often *scaled* by the square root of the vector dimensionality:

$$f(\mathbf{q}, \mathbf{k}) = \frac{1}{\sqrt{D}} \langle \mathbf{q}, \mathbf{k} \rangle \quad (5.130)$$

The projection function and soft and hard attention. The projection function used to transform the un-normalized attention scores is a crucial component of the transformer model. By far the most commonly used projection function is again the softmax. In this case, the attention function is referred to as soft attention.

Definition 5.2.5 (Soft attention). *The **soft attention** is computed with the projection function $\mathbf{f}_{\Delta^{D-1}} = \text{softmax}$.*

However, the softmax again makes the models difficult to analyze. In our voyage to theoretically understand transformer-based language models, we will, therefore, again make specific (less frequently used) modeling choices, particularly in the case of the projection function.

Indeed, to be able to derive any interesting expressivity results (see §5.3), we jump to the other side of the spectrum and define hard attention. Simply put, instead of spreading the attention across all values like softmax, hard attention puts all the mass on the element whose key maximizes the scoring function f . One way to arrive at it from the definition of soft attention is by sending the temperature τ in the definition of the softmax function (cf. Definition 3.1.10) to 0. Recall that that results in the output vector representing a uniform distribution over the elements that maximize the input vector. This is known as *averaging* hard attention.

Definition 5.2.6 (Averaging hard attention). *The **averaging hard attention** is an attention mechanism with the projection function $\mathbf{f}_{\Delta^{D-1}} = \text{hardmax}_{\text{avg}}$, where $\text{hardmax}_{\text{avg}}$ is defined as:*

$$\text{hardmax}_{\text{avg}}(\mathbf{x})_d \stackrel{\text{def}}{=} \begin{cases} \frac{1}{r} & \text{if } d \in \text{argmax}(\mathbf{x}) \\ 0 & \text{otherwise} \end{cases}, \text{ for } d = 1, \dots, D \quad (5.131)$$

where $\mathbf{x} \in \mathbb{R}^D$ and $r = |\text{argmax}(x)|$ is the cardinality of the *argmax* set over \mathbf{x} .

Interestingly, there is another form of hard attention that results in a model with a different expressive capacity: the unique hard attention. The difference lies exactly in how it handles *ties* in the elements which maximize the scoring function. Unique hard attention chooses only *a single* element of those that maximize the score: it can be chosen randomly or deterministically (e.g., always the first one).

Definition 5.2.7 (Unique hard attention). *The **unique hard attention** is an attention mechanism with the projection function $\mathbf{f}_{\Delta^{D-1}} = \text{hardmax}_{\text{uni}}$, where $\text{hardmax}_{\text{uni}}$ is defined as follows. For $\mathbf{x} \in \mathbb{R}^D$, sample $\hat{d} \sim \text{Unif}(\text{argmax}(\mathbf{x}))$ or choose some $\hat{d} \in \text{argmax}(\mathbf{x})$ deterministically. Then*

$$\text{hardmax}_{\text{uni}}(\mathbf{x})_d \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } d = \hat{d} \\ 0 & \text{otherwise} \end{cases}, \text{ for } d = 1, \dots, D. \quad (5.132)$$

³⁴While the fact that the transformed scores sum to one invites their interpretation as probabilities, this is not their central role. Rather, the weights are simply used to define a convex combination of the values.

While the difference between unique and averaging hard attention might seem subtle and marginal, it actually results in a large difference in the expressivity of transformer-based language models as we discuss in §5.3. While we will investigate this in a lot more detail there, we just mention that the intuition behind the expressive difference is relatively straightforward: while the *keys* maximizing the un-normalized scores might be the same (even though they necessarily don't have to be if f is not injective), the *values* (whose content is decoupled from the keys) that those keys index might not be—and in some cases, all those different values might be relevant for the task at hand. Unique hard attention always allows us to only “lookup” a *single* value associated with those keys, no matter how “different” and relevant all of those are. It also does not allow the attention mechanism to sum over (“summarize”) across all the elements that maximize attention. This is a very limiting characteristic, as many of the expressivity results that we will see later rely on summing over all the elements that maximize the attention scores.

Transformer Blocks

We have, so far, described the “low-level” details of how the attention function is computed. We now combine the computations performed into larger blocks, showing how they are used to compute the string-augmented encodings of the individual symbols. In particular, we have to connect the concepts of queries, keys, and values to the symbols and their (initial, static) encodings. Intuitively, this is done by transforming the static encodings of those symbols using specific functions implemented through the attention mechanism and using the transformed encodings as queries, keys, and values, as we describe below.

We first abstract the attention mechanism from Definition 5.2.4 a bit. With this, we will, in a few steps, arrive at exactly how the hidden states \mathbf{h}_t or the contextual encodings are computed in a transformer. At the core of this computation lies a repeated application of the same sequence of operations, which, as we will see, augment the “current version” of the contextual encodings with the current information from the preceding information. We call a single sequence of operations a transformer layer.

Definition 5.2.8 (Transformer layer). *Let Q , K , V , and O be parametrized functions from \mathbb{R}^D to \mathbb{R}^D .*

*A **transformer layer** is a function $T: \mathbb{R}^{T \times D} \rightarrow \mathbb{R}^{T \times D}$ that takes as input sequence of vectors $\mathbf{X} = (\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_T^\top)$ and returns $\mathbf{Z} = (\mathbf{z}_1^\top, \mathbf{z}_2^\top, \dots, \mathbf{z}_T^\top) \in \mathbb{R}^{T \times D}$ according to the following steps:*

$$\mathbf{a}_t = \underbrace{\text{Att}(Q(\mathbf{x}_t), K(\mathbf{X}_t), V(\mathbf{X}_t))}_{\text{Definition 5.2.4}} + \mathbf{x}_t \quad (5.133)$$

$$\mathbf{z}_t = O(\mathbf{a}_t) + \mathbf{a}_t \quad (5.134)$$

for $t = 1, \dots, T$, so that $T(\mathbf{X}) \stackrel{\text{def}}{=} \mathbf{Z} = (\mathbf{z}_1^\top, \mathbf{z}_2^\top, \dots, \mathbf{z}_T^\top) \in \mathbb{R}^{T \times D}$.

While we defined the transformer layer on a general matrix (with T columns), note that these T vectors will refer to the (current) symbol encodings of the symbols in the string up to the T^{th} symbol, i.e., $\mathbf{y}_{\leq T}$.

What do these quantities correspond to? Eqs. (5.133) and (5.134) outline a two-step process of computing the outputs of a *single* transformer layer: $\mathbf{X} = (\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_T^\top)$ represents the input to the layer, which T transforms into the output sequence $\mathbf{Z} = (\mathbf{z}_1^\top, \mathbf{z}_2^\top, \dots, \mathbf{z}_T^\top)$. Before being fed into the attention mechanism, the inputs \mathbf{X} are first *transformed* into the quantities required by the attention mechanism: the query \mathbf{q}_t (a single one for each \mathbf{x}_t), the matrix of keys \mathbf{K}_t , and the matrix of values \mathbf{V}_t —all of these are, therefore, *transformations* of the input sequence of vectors. The transformations Q , K , and V determine how these inputs are transformed into the (interpretable) quantities required by the attention mechanism.

The individual \mathbf{a}_t represent the “intermediate” results of the computation—the results of applying the actual attention mechanism (with a slight modification) from Definition 5.2.4 onto the produced values of the query, the keys, and the values.

The modification mentioned refers to the addition of the inputs \mathbf{x}_t to the output of the attention mechanism in Eq. (5.133). This mechanism is known as adding **residual connections** to the model. First introduced by He et al. (2016) in the context of deep convolutional neural network-based architectures, residual connections are now a common feature in many state-of-the-art deep learning architectures. Note that their use is mostly motivated by empirically better performance—this is often attributed to the fact that, intuitively, residual connections allow gradients (i.e., learning signals) to flow through the network through a more direct route rather than all layers that can “squish” the signal similarly to the Elman RNN case (in that sense, they help

mitigate the vanishing gradient issue). However, as we will see later in our analysis, residual connections will also play an important role in determining the theoretical properties of transformers, particularly their computational expressive power. The same mechanism is applied in the second step of the transformer layer, where the intermediate results \mathbf{a}_t are transformed by the output transformation O into the final outputs of the layer.

In the simplest case, you can imagine the inputs \mathbf{X} to be the initial static embeddings of the symbols. The application of the transformer layer, in this case, therefore, “selectively” (determined by the attention mechanism) augments the static embeddings with the information from the preceding context. However, as we will see shortly, a transformer model will apply *multiple* transformer blocks to the input sequence and thus transform it in multiple steps, analogously to how layers are composed in a regular feed-forward neural network. In that sense, the inputs to the transformer blocks will refer to general intermediate representations produced from the initial static embeddings after some number of applications of transformer layers.

Lastly, let us consider how the current symbol representations \mathbf{X} are transformed into the queries, keys, and values using Q , K , and V ? The original formulation (Vaswani et al., 2017) and all standard implementations of the transformer architecture use one of the simplest possible mappings: a linear transformation implemented by matrix multiplication. On the other hand, the final output of the transformer, computed with output mapping O , is usually implemented by a multi-layer perceptron.

The transformer layer puts the attention mechanism into a functional block that describes how a sequence of current symbol representations is transformed in a single step into a sequence of representations augmented with the information in the current set of values. However, we are not done abstracting yet! As mentioned, this process can be applied arbitrarily many times, resulting in “deep” encodings of individual symbols, which contain information from the preceding symbols in the string computed in a composite way. This also answers the question of how the augmented symbol representations used in the language modeling formulation are computed from the initial symbol representations: they are the result of multiple applications of the transformer layer to the matrix of initial symbol representations. That is: multiple transformer layer layers are stacked on top of one another so that the output of one layer becomes the input of the next.

We now have all the building blocks to define the full transformer architecture, which computes the encodings of the string prefixes (and thus the hidden states) in Definition 5.2.3.

Definition 5.2.9 (Transformer). *For $L \in \mathbb{N}$, we define a L -layer **transformer** model as a D -dimensional transformer sequence model over an alphabet Σ where the hidden state $\mathbf{h}_t \stackrel{\text{def}}{=} \text{enc}_{\mathcal{T}}(y_1 \dots y_t) = \text{enc}_{\mathcal{T}}(\mathbf{y})$ is computed as follows.*

$$\mathbf{X}^1 \stackrel{\text{def}}{=} (\mathbf{e}'(y_0), \mathbf{e}'(y_1), \dots, \mathbf{e}'(y_t)) \quad (5.135)$$

$$\mathbf{Z}^\ell = \mathcal{T}_\ell(\mathbf{X}^\ell) \text{ for } 1 \leq \ell < L \quad (5.136)$$

$$\mathbf{X}^{\ell+1} = \mathbf{Z}^\ell \text{ for } 1 \leq \ell < L \quad (5.137)$$

$$\mathbf{h}_t = F(\mathbf{z}_t^L) \quad (5.138)$$

\mathcal{T}_ℓ for $\ell = 1, \dots, L$ represent L different transformer layers with decoupled parameter (cf. Definition 5.2.8). $F: \mathbb{R}^D \rightarrow \mathbb{R}^D$ is a transformation function applied to the contextual encoding of the last symbol in the last (L^{th}) layer and $\mathbf{e}': \Sigma \rightarrow \mathbb{R}^D$ is a symbol representation function computing the initial representations of the symbols passed to the first layer of the transformer.³⁵

With this, the transformer model now fully specifies how to compute the representations required for the representation-based locally normalized sequence models from Definition 5.2.3—the representation function $\text{enc}_{\mathcal{T}}$ is the composition of L transformer layers applied to the sequence of static encodings, followed by a final transformation F .

Making Attention Work Fast Over Entire Strings

Notice that, in the formulations so far, we always presented computations of the attention mechanism for individual queries \mathbf{q}_i . This corresponds to the computation of the new version of the representation of a

³⁵The symbol representation function \mathbf{e}' is often also implemented as a linear transformation of the one-hot representations (cf. Definition 5.1.5) of symbols, i.e., it is simply a table-lookup.

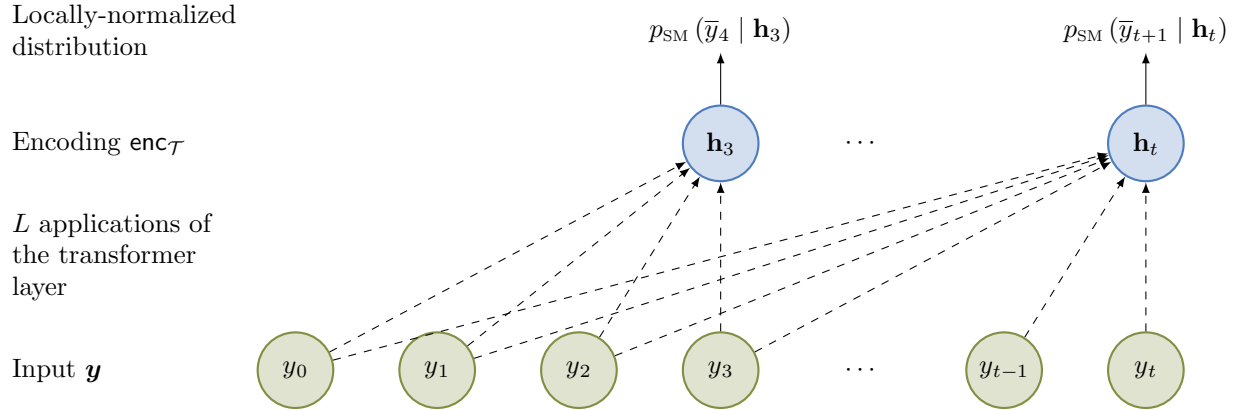


Figure 5.20: $\text{enc}_T(\mathbf{y}_{\leq t})$ is a function of the symbols y_1, \dots, y_t computed with multiple applications of the transformer block. Here, the dashed lines illustrate the dependencies of the outputs \mathbf{h}_t on the initial static encoding of the symbols y denoted by green nodes. Naïvely, \mathbf{h}_3 and \mathbf{h}_t could be computed by independently applying the attention mechanism from Definition 5.2.4. However, as we describe in the text, while the applications of the attention mechanism do not *share* computations, they can be written concisely together.

single symbol in the string, with the keys and values representing the preceding symbols (including the symbol itself). This is illustrated in Fig. 5.20. This could of course be applied $|\mathbf{y}|$ -times to compute the representations of all $|\mathbf{y}|$ symbols in a single transformer block. However, this would unnecessarily re-compute the keys and values of the symbols multiple times—and, as we motivated at the beginning, speedups were one of the main reasons to talk about transformers in the first place. We can now show how the attention mechanism can be conveniently applied to entire strings at once. Specifically, we focus on the case where the attention scoring function f is implemented as a dot-product.³⁶

What does the attention mechanism from Definition 5.2.4 do in this case? Given a query \mathbf{q}_t and a matrix of key values $\mathbf{K} = (\mathbf{k}_1^\top, \dots, \mathbf{k}_t^\top) \in \mathbb{R}^{t \times D}$, the scoring function simply computes³⁷

$$u_j = f(\mathbf{q}_t, \mathbf{k}_j) = \mathbf{q}_t^\top \mathbf{k}_j.$$

Notice that, in this case, the vector $\mathbf{u} = (u_1, \dots, u_t)$ of *unnormalized* attention weights can simply be computed as a single matrix-vector product

$$\mathbf{u} = \mathbf{q}_t^\top \mathbf{K}^\top.$$

Furthermore, with this, attention can be easily extended to consider many queries in parallel by stacking multiple queries into a matrix $\mathbf{Q} \stackrel{\text{def}}{=} (\mathbf{q}_1^\top, \mathbf{q}_2^\top, \dots, \mathbf{q}_t^\top)$, as we detail now.³⁸ Consider now the product

$$\mathbf{U} = \mathbf{Q} \mathbf{K}^\top.$$

Each entry of the resulting matrix U_{ij} is exactly the dot-product between the query \mathbf{q}_i and the key \mathbf{k}_j ! The *rows* of \mathbf{U} then contain the unnormalized score vectors \mathbf{u}_i from the definition of the attention mechanism. This means that if we now apply the normalization function $\mathbf{f}_{\Delta^{D-1}}$ *row-wise* (such that the sums of the elements in each row equal 1), we end up with exactly the required normalized scores required for combining the values from the value matrix. With some abuse of notation, we will simply write that as

$$\mathbf{S} \stackrel{\text{def}}{=} (\mathbf{s}_1^\top, \dots, \mathbf{s}_t^\top) \stackrel{\text{def}}{=} \mathbf{f}_{\Delta^{D-1}}(\mathbf{U}) = \mathbf{f}_{\Delta^{D-1}}(\mathbf{Q} \mathbf{K}^\top). \quad (5.139)$$

³⁶For conciseness, we will ignore the scaling factor, which could easily be added.

³⁷We switch notation from $\langle \mathbf{q}_t, \mathbf{k}_j \rangle$ to $\mathbf{q}_t^\top \mathbf{k}_j$ to make the connection to matrix multiplication later clearer.

³⁸Note that, for easier presentation, we make a slight departure from the original definition of the attention mechanism, where the result of the attention mechanism for query t only depended on the keys and values $j \leq t$. For the rest of the paragraph, we assume that a query \mathbf{q}_i with $i < t$ can consider keys and values \mathbf{k}_j and \mathbf{v}_j with $j > i$, which, in the interpretation of attention applied to strings, would mean that the symbols can “look ahead” in the string to their right. This will be addressed shortly with *masking*.

The rows of $\mathbf{f}_{\Delta^{D-1}}(\mathbf{A})$, therefore, represent the normalized attention weights. This brings us to the final step of the matrix-multiplication-based attention mechanism: Combining the values based on the computed attention weights. Again, this can be performed by a single matrix multiplication. Notice that the *value* vectors are *the same* for all queries—they are simply combined with different (attention) weights based on the query. Right-multiplying the *transposed* values matrix $\mathbf{V} = (\mathbf{v}_1^\top, \dots, \mathbf{v}_t^\top)$ with \mathbf{S} , therefore, perform the convex combination of the value vector $\mathbf{v}_1^\top, \dots, \mathbf{v}_t^\top$ such that

$$\mathbf{a}_i = \mathbf{s}_i \mathbf{V}^\top = \mathbf{S}_{i,:} \mathbf{V}^\top \quad (5.140)$$

and thus

$$\mathbf{A} \stackrel{\text{def}}{=} (\mathbf{a}_1, \dots, \mathbf{a}_t) = \mathbf{S} \mathbf{V}^\top. \quad (5.141)$$

Altogether, this means that, given a sequence of (contextual) symbol encodings \mathbf{X} , we can compute the attention values (i.e., the output of the attention mechanism) of *all* queries (i.e., for all string in the string) with a single matrix multiplication, as long as the scoring function is the (scaled) dot-product. We refer to this version of attention as an attention block, which, intuitively, simply replaces the *element-wise* definition of the attention mechanism from Definition 5.2.4 with a more efficient (and concise) definition through matrix multiplications.³⁹

Definition 5.2.10 (Attention Block). *Let Q , K , and V be parametrized functions from $\mathbb{R}^{T \times D}$ to $\mathbb{R}^{T \times D}$ and $\mathbf{X} \in \mathbb{R}^{T \times D}$ the matrix of input encodings. An **attention block** is the function $\mathbf{A}: \mathbb{R}^{T \times D} \rightarrow \mathbb{R}^{T \times D}$ defined as*

$$\mathbf{A}(\mathbf{X}) = \mathbf{f}_{\Delta^{D-1}} \left(Q(\mathbf{X}) K(\mathbf{X})^\top \right) V(\mathbf{X}) \quad (5.142)$$

Further, we define the **attention matrix** as the square matrix $\mathbf{U} \stackrel{\text{def}}{=} Q(\mathbf{X}) K(\mathbf{X})^\top \in \mathbb{R}^{T \times T}$.

As mentioned, the functions $Q(\cdot)$, $K(\cdot)$, and $V(\cdot)$ are usually implemented as a linear transformation via matrix multiplication using weight matrices \mathbf{W}_Q , \mathbf{W}_K , and \mathbf{W}_V . This means that the query matrix \mathbf{Q} can be computed as $\mathbf{Q} = \mathbf{X} \mathbf{W}_Q$, where $\mathbf{W}_Q \in \mathbb{R}^{D \times D}$ is a matrix of learnable parameters. Since the attention block uses the same input matrix \mathbf{X} to encode queries, keys, and values, it is usually called **self-attention**.

Confining attention to preceding symbols in the string. We now address the departure from the original definition of the attention mechanism in which the query \mathbf{q}_t was only allowed to consider the keys and values \mathbf{k}_j and \mathbf{v}_j with $j \leq t$. Notice that, in general, the version of attention of Eq. (5.142) allows each symbol to attend to *any* symbol in the string, even those in later positions in the string. Note that there is nothing inherently *wrong* with that—the contextual symbol encodings could, in principle, depend on the information from the *entire* string. In fact, this is very common in the so-called *masked language modeling*, which, importantly, despite its name, does not define language models in our sense of the word. A very commonly used family of masked models is the BERT family of models.⁴⁰ However, in the case of locally normalized language models, “looking ahead” obviously violates the autoregressive structure of language modeling, i.e., violates our assumption that the context at time t forms $\mathbf{y}_{<t}$. This is illustrated in Fig. 5.21. To recover the autoregressive nature of the language model, we, therefore, posthoc modify Eq. (5.142) to allow each symbol to attend only to *itself* and to *preceding* symbols, while still being able to implement it using matrix multiplication. We do that by adding a *mask* to zero out the unwanted elements of \mathbf{U} .

Definition 5.2.11 (Masked Attention Block). *Let $Q(\cdot)$, $K(\cdot)$, and $V(\cdot)$ be parametrized functions from $\mathbb{R}^{T \times D}$ to $\mathbb{R}^{T \times D}$. A **masked attention block** is a function $\mathbf{A}(\mathbf{X}, \mathbf{M}): \mathbb{R}^{T \times D} \times \mathbb{R}^{T \times D} \rightarrow \mathbb{R}^{T \times D}$ defined as*

$$\mathbf{A}(\mathbf{X}, \mathbf{M}) = \text{softmax}(Q(\mathbf{X}) K(\mathbf{X})^\top \odot \mathbf{M}) V(\mathbf{X}) \quad (5.143)$$

where \odot is the element-wise product between matrices, and $\mathbf{M} \in \mathbb{R}^{\ell \times \ell}$, the **masking matrix**, is constructed as follows.

$$M_{i,j} = \begin{cases} 1 & \text{if } i \leq j \\ -\infty & \text{otherwise} \end{cases} \quad \text{for } 0 \leq i, j < T \quad (5.144)$$

³⁹Again, with the caveat that the attention weights are not confined to the preceding symbols but to all symbols in the string.

⁴⁰In a very unfortunate, but also understandable, turn of events, we mention two completely different notions of masking in a single paragraph. Importantly, the masking in *masked language models* (which, again, are not language models in our strict definition) has nothing to do with the “causal” masking relevant for autoregressive language modeling which we introduce in this section.

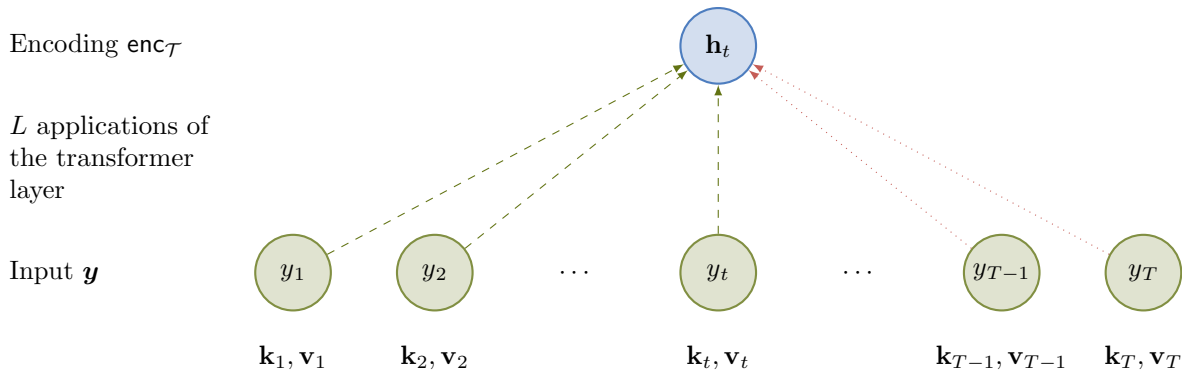


Figure 5.21: In the context of language modeling, the attention mechanism is allowed to consider the symbols y_j and their keys/values for $j \leq t$ (the **green dashed lines**) when computing the contextual encoding of the symbol y_t . In Definition 5.2.4 this is enforced by the definition of the matrices \mathbf{K}_t and \mathbf{V}_t . However, in the attention block formulation of the attention mechanism from Definition 5.2.10, since the matrices \mathbf{K} and \mathbf{V} contain the values corresponding to the entire string \mathbf{y} , the query \mathbf{q}_t could, in principle, index into the values corresponding to the symbols y_{t+1}, \dots, y_T (the **red dotted lines**). Masking prevents that by enforcing the attention weights a_{t+1}, \dots, a_T to be 0. In this sense, it removes the **red dotted lines**.

This implements a very easy “fix” to the looking-ahead problem by simply putting the normalized attention scores of the “illegal” elements to 0. In general, the exact value of the elements $M_{i,j}$ with $i > j$ of course depends on the projection function $\mathbf{f}_{\Delta^{D-1}}$ —for simplicity, we only define \mathbf{M} for the case where $\mathbf{f}_{\Delta^{D-1}} = \text{softmax}$.

Bits and Bobs of the Transformer Architecture: Positional Encodings, Multiple Heads, and Layer Normalization

Let us now take a step back and consider what the transformer model introduced so far does abstractly. A transformer takes as input a string $\mathbf{y} \in \Sigma^*$, computes the initial symbol embeddings $\mathbf{X}^1 = \mathbf{X}$, and transforms those through a sequence of L applications of the transformer layer (cf. Definition 5.2.9). This results in the final augmented (contextual) symbol representations $\mathbf{h}_t = \text{enc}_{\mathcal{T}}(\mathbf{y}_{\leq t})$, which are then used to compute the conditional probabilities in the representation-based locally normalized language model defined by the transformer (cf. Definition 5.2.3), as illustrated on top of Fig. 5.20. In this subsection, which will finish off our formal definition of the architecture, we introduce the last three components often connected closely to the transformer model: Symbol positional encodings, multi-head attention, and layer normalization.

Adding positional information into the transformer architecture. There is an important omission we still have not addressed when talking about transformers: How does the model incorporate any notion of *word order* into the contextual representations of symbols or the encodings of the context \mathbf{h}_t ? The motivation is very clear: The meaning of a sentence depends on the word order. The meaning of “A dog bit a man.” is not the same as “A man bit a dog.”. This is one of the reasons why simple “sentence encoding” functions such as bag-of-words, which simply represent sentences with the number of individual words they contain, do not work well. A careful reader might have noticed that at no point in our discussion about transformers and the attention mechanism did we say anything about the positions of the words. Importantly, we did not talk about word positions in the case of RNNs either. However, the sequential and incremental processing nature of RNNs makes it easy to “manually” keep track of the position of the current symbol of the string y_t , to the extent that the RNN variant is capable of “counting” (cf. §5.1.7). However, all operations composing the transformer model are *position-agnostic*: The convex combination of the value vectors \mathbf{V} will be the same, no matter the permutation of the vectors (if we, of course, accordingly permute the keys). The keys also cannot contain any positional information, since they are computed from position-agnostic static embeddings and a transformation function K which does not depend on the position.

All that is to say that, to be able to take into account word order in a transformer, we have to *explicitly*

provide the positional information to the model. The simplest way to do this is to *augment* the static symbol encodings in the first transformer layer with *positional encodings* in the form of vectors which can be added to or concatenated to the static encodings of symbols (Vaswani et al., 2017).⁴¹

Definition 5.2.12 (Positional encoding). A **positional encoding** is a function $\mathbf{f}_{\text{pos}}: \mathbb{N} \rightarrow \mathbb{R}^D$.

This is a very simple definition: A positional encoding simply assigns a position in a string a vector representation. A trivial example would be $\mathbf{f}_{\text{pos}}(t) = (t)$. This allows us to define a position-augmented symbol representation function.

Definition 5.2.13 (Position-augmented representation function). Let $\mathbf{e}': \bar{\Sigma} \rightarrow \mathbb{R}^D$ be a symbol representation function and $\mathbf{f}_{\text{pos}}: \mathbb{N} \rightarrow \mathbb{R}^D$ a positional encoding. A **position-augmented** representation function of a symbol y_t in a string \mathbf{y} is the representation function $\mathbf{e}'_{\text{pos}}: \bar{\Sigma} \rightarrow \mathbb{R}^D$ defined as

$$\mathbf{e}'_{\text{pos}}(y_t) \stackrel{\text{def}}{=} \mathbf{e}'(y_t) + \mathbf{f}_{\text{pos}}(t). \quad (5.145)$$

To make the positional information available to the transformer model, we now simply pass the position-augmented “static” symbol encodings \mathbf{X}_{pos} to the model instead of the original ones \mathbf{X} . Apart from that, the transformer model can remain unaltered, and function simply as defined above, taking into account the positional information now included in its inputs. Importantly, the intuitive notion of the importance of positional encodings for understanding natural language also transfers to the computational power of the model: Transformers as introduced in this section *without* positional information are strictly less powerful than those with positional information (Pérez et al., 2021). Again, this intuitively makes sense: Without positional information, a transformer model could not even recognize the simple (unweighted) regular language

$$L = \{ab^n \mid n \in \mathbb{N}\}$$

since it would have no way of knowing, provided that a is in a given string \mathbf{y} , whether it appears in the first position or in any other position in the string.

Multiple heads. Importantly, the transformer introduced so far computes a single set of contextual representations—one for every input symbol (at every layer of the transformer). However, we can easily extend the model to compute *multiple* contextual representations for each symbol. This is done using the so-called **multi-head attention**, where a single attention block is called an **attention head**. This increases the representation space of the individual symbols and thus enables the model to capture more information about the symbols and the sentence. The interpretation of computing *multiple* representations (one for each head) independently also invites the interpretations that each of the heads “focuses” on a separate aspect of the text. To be able to use the outputs of multi-head attention as inputs to the next block again, the outputs of the different attention heads are then concatenated and then projected down to the output size of a single attention block using an additional transformation.

Definition 5.2.14 (Multi-Head Attention Block). Let $H \in \mathbb{N}$ be the number of attention heads, $Q_h(\cdot)$, $K_h(\cdot)$, and $V_h(\cdot)$ be parametrized functions from $\mathbb{R}^{T \times D}$ to $\mathbb{R}^{T \times D}$ for $0 \leq h \leq H$, and $\mathbf{f}_H: \mathbb{R}^{T \cdot H \times D} \rightarrow \mathbb{R}^{T \times D}$ be a parametrized function.

A **multi-head attention block** is a function $\text{MH-A}(\mathbf{X}): \mathbb{R}^{T \times D} \rightarrow \mathbb{R}^{T \times D}$ defined as

$$\text{MH-A}(\mathbf{X}) = \mathbf{f}_H(\text{concat}_{0 \leq h < H} \left(\text{softmax} \left(Q_h(\mathbf{X}) K_h(\mathbf{X})^\top \right) V_h(\mathbf{X}) \right)) \quad (5.146)$$

While multi-head attention is mostly motivated by empirically better performance (and the intuitive motivation of being able to separately focus on different notions of similarity), it will have some implications on the computational power of the model as well. As we will see shortly, having multiple heads makes it very easy to reason about how a transformer model can simulate an n -gram model.⁴²

⁴¹For simplicity, we assume the positional encodings are added to the static ones; notice that by dividing the D -dimensional vectors into two components, one responsible for the static encodings and one for the positional ones (where the positional encoding component is zeroed out in the static encoding and vice-versa), one can easily implement “concatenation” of the two representations using only addition.

⁴²This does not, however, mean that multiple heads are *required* for recognizing n -gram models. As we will see, under some caveats, single-head transformers are Turing complete.

Layer normalization. As a final component of a transformer, we mention layer normalization. Layer normalization, similar to the use of residual connections, represents a common “trick” in the deep learning space for ensuring more stable and reliable gradient-based learning—as such, it is not limited to transformers. Formally, we can define layer normalization as follows (Ba et al., 2016).

Definition 5.2.15. Let $\mathbf{x}, \gamma, \beta \in \mathbb{R}^D$, and $\epsilon > 0$. The **layer normalization** function $\text{LN}: \mathbb{R}^D \rightarrow \mathbb{R}^D$ is defined as

$$\text{LN}(\mathbf{x}; \gamma, \beta) \stackrel{\text{def}}{=} \frac{\mathbf{x} - \bar{\mathbf{x}}}{\sqrt{\sigma^2(\mathbf{x}) + \epsilon}} \odot \gamma + \beta, \quad (5.147)$$

where $\bar{\mathbf{x}}$ refers to the mean of the vector \mathbf{x} (and is subtracted from all elements of \mathbf{x} in the formulation above) and $\sigma^2(\mathbf{x})$ refers to the variance of elements of \mathbf{x} . ϵ is added in the denominator to ensure stability if $\sigma^2(\mathbf{x}) \ll 1$.

Intuitively, the application of the layer normalization function ensures that the mean of the vector \mathbf{x} is (approximately) β and its variance is controlled by γ (after being standardized by dividing by the standard deviation of \mathbf{x}). Most commonly, we simply set $\gamma = \mathbf{1} \in \mathbb{R}^D$ and $\beta = \mathbf{0} \in \mathbb{R}^D$.

Layer normalization is most commonly applied to the *output of the transformer layer* (on every layer), i.e., to \mathbf{z}_i in Eq. (5.134). The full output of the transformer layer is therefore computed as

$$\mathbf{z}_i \stackrel{\text{def}}{=} \text{LN}(O(\mathbf{a}_i) + \mathbf{a}_i; \gamma, \beta). \quad (5.148)$$

Interestingly, although we mentioned that layer normalization is mostly motivated by the stability it brings to training and with it better performance, it does, just like multi-headed attention, seem to contribute to the computational expressivity of transformer models. As we will see in §5.3, layer normalization allows for a simple fix that solves one of the best known formal limitations of the transformer architecture (again, under some assumptions) (Hahn, 2020; Chiang and Cholak, 2022).

Connecting the Formal Definition Back to our Desiderata

This brings us to the end of the formal definition of the transformer architecture. As we saw, a transformer has many more moving parts than an RNN. Is the additional complexity warranted? Let us now return back to the informal motivations or desiderata that we laid out in §5.2.1 and see how the components we defined in this section come through and ensure transformers fit them. First of all, the transformer layers clearly store the representations of all symbols at all times—they are all needed to produce the query, key, and value matrices required by the attention mechanism. As mentioned above, this allows us to easily store information about the entire string in a convenient and “accessible” format without having to compress it into a single hidden state. Furthermore, the fact that the contextual representations $\mathbf{Z}^{(\ell+1)}_t$ are computed from the representations $\mathbf{x}_1^\ell, \dots, \mathbf{x}_t^\ell$ directly at every time step, with no direct dependence between the different $\mathbf{z}^{(\ell+1)}_i$, means that these computations can easily be parallelized. More concretely, in the most common implementations of the transformer components, most operations take the form of matrix multiplications, which makes the computation and parallelization that much more efficient.⁴³ Again, note that here, we are only interested in parallelizing the processing of *entire* strings, as for example given in a training corpus. As discussed in §5.1.6, there is an aspect of language modeling that is inherently sequential and even heavily parallelizable architectures such as the transformer cannot overcome: Generating strings one symbol at time. While generating strings, even a transformer model will have to generate symbols one at a time and, therefore, recompute (parts of) the encoding $\text{enc}_T(\mathbf{y}_{\leq t})$ anew at every time step to generate the next symbol. The advantages of parallelizability, therefore, come only at training time—however, given the vast corpora used for training today’s models, this makes a crucial difference in the applicability of the architecture over recurrent ones.

Altogether, this means that transformers do, indeed, achieve the desiderata from our informal motivation! This concludes our formal definition of transformers. We move to analyze their theoretical properties.

⁴³Note that, there is, of course, some sense of recurrence in the transformer—the composition of the transformer layers, which are stacked on top of each other, and of course require sequential computation. However, crucially, the number of layers does not depend on the *length of the string*—the number of sequential steps required to process a string, therefore, does not depend on its length, which is what we wanted to achieve.

5.2.3 Tightness of Transformer-based Language Models

Having introduced transformers formally, we can start investigating their formal properties. As we did for RNN LMs, we first consider their tightness. Specifically, in this subsection, we show that all *soft attention-based* transformer language models are tight. Key to our proof of the tightness of transformer language models, as well as the tightness of various other neural architectures, is the following basic fact in topology.

Theorem 5.2.1. *Let \mathcal{X} be a compact topological space and \mathcal{Y} be any topological space. If $f: \mathcal{X} \rightarrow \mathcal{Y}$ is continuous, then $f(\mathcal{X}) \subseteq \mathcal{Y}$ is also compact.*

Proof. Let $\{\mathcal{U}_s\}_{s \in \mathcal{A}}$ be any open cover of $f(\mathcal{X})$. By continuity, $f^{-1}(\mathcal{U}_\alpha) \subset \mathcal{X}$ is open for any $\alpha \in \mathcal{A}$, and hence $\{f^{-1}(\mathcal{U}_s)\}_{s \in \mathcal{A}}$ is also an open cover of \mathcal{X} . By the compactness of \mathcal{X} , there is a finite sub-cover $\{f^{-1}(\mathcal{U}_{\alpha_n})\}_{n=1}^N$, in which case $\{\mathcal{U}_{\alpha_n}\}_{n=1}^N$ forms a finite sub-cover for $f(\mathcal{X})$. ■

We now further mathematically abstract transformers as a function on vector tuples,⁴⁴ $\mathbf{f}_{\text{Att}}: (\mathbb{R}^D)^+ \rightarrow (\mathbb{R}^D)^+$, that is *length-preserving* in the sense that $\mathbf{f}_{\text{Att}}(\mathbb{R}^{t \times D}) \subseteq (\mathbb{R}^{t \times D})$ for all $t > 0$. Intuitively, this definition is saying that \mathbf{f}_{Att} is a function that maps a nonempty vector tuple $\{\mathbf{v}_j\}_{j=1}^t$ to another vector tuple $\{\mathbf{h}_j\}_{j=1}^t$ of the same length,

$$\mathbf{f}_{\text{Att}}(\mathbf{v}_1, \dots, \mathbf{v}_t) = (\mathbf{h}_1, \dots, \mathbf{h}_t) \in \mathbb{R}^{t \times D}, \quad (5.149)$$

where $\mathbf{v}_j = \mathbf{e}'(y_j) \in \mathbb{R}^D$ are the initial representations of the input symbols y_j . In particular, we can take the function $\mathbf{f}_{\text{Att}}: (\mathbb{R}^D)^+ \rightarrow (\mathbb{R}^D)^+$ to be the function defined by a stack of transformer layers, i.e., an attention block. This setup will help us state the following.

Lemma 5.2.1. *Let $\mathbf{f}_{\text{Att}}: (\mathbb{R}^D)^+ \rightarrow (\mathbb{R}^D)^+$ be the function defined by a L transformer layers with continuous functions Q, K, V , and O . Given a compact set $\mathcal{K} \subset \mathbb{R}^D$. Then, there exists a compact set $\mathcal{K}' \subset \mathbb{R}^D$ such that for every $t \in \mathbb{Z}_{>0}$,*

$$\mathbf{f}_{\text{Att}}(\mathcal{K}^t) \subseteq (\mathcal{K}')^t. \quad (5.150)$$

Note. We make use of the following notations in the proof below: $B_r(\mathbf{z}) = \{\mathbf{v} \in \mathbb{R}^D : \text{dist}(\mathbf{z}, \mathbf{v}) < r\}$ denotes the open ball centered at \mathbf{z} with radius r ; $\overline{\mathcal{A}}$ denotes the closure of set \mathcal{A} .

Proof. Let $\mathcal{K}_0 = \mathcal{K}$. In an autoregressive transformer, each of the L layers consists of two blocks: a self-attention block and a feedforward block. We will use induction on the $2L$ blocks to build up compact sets $\mathcal{K}_1, \mathcal{K}_2, \dots, \mathcal{K}_{2L}$ that contain the output vectors of these respective blocks, and then take $\mathcal{K}' = \mathcal{K}_{2L}$.

The self-attention block is a function on $(\mathbb{R}^D)^+ \rightarrow (\mathbb{R}^D)^+$. So, let $t \in \mathbb{Z}_{>0}$ be arbitrary and consider any sequence of input vectors $(\mathbf{v}_1, \dots, \mathbf{v}_t)$ such that for all i , $\mathbf{v}_i \in \mathcal{K}_0$. Denote the output vectors of the attention block with $(\mathbf{v}'_1, \dots, \mathbf{v}'_t)$. By definition of attention, each output vector $\mathbf{v}'_j = \sum_{i=1}^t s_i^{(j)} \mathbf{v}_i$ where $\mathbf{s}^{(j)} \in \Delta^{t-1}$ are the attention weight vectors obtained through the softmax function. Compact sets in \mathbb{R}^D are bounded (by the Heine–Borel theorem), and hence there exists $M > 0$ such that $\mathcal{K}_0 \subseteq \overline{B_M(0)}$. Noting that the norm function $\|\cdot\|$ on \mathbb{R}^D is convex, we have the following

$$\|\mathbf{v}'_j\| = \left\| \sum_{i=1}^t s_i^{(j)} \mathbf{v}_i \right\| \quad (5.151a)$$

$$\leq \sum_{i=1}^t s_i^{(j)} \|\mathbf{v}_i\| \quad (*)$$

$$\leq \sum_{i=1}^t s_i^{(j)} M = M \quad (5.151b)$$

where $(*)$ results from Jensen's inequality. Eq. (5.151b) shows that each of the output vectors \mathbf{v}'_j lies in $\overline{B_M(0)}$ which is compact. Hence, setting $\mathcal{K}_1 = \overline{B_M(0)}$, we have shown that, for any $t \in \mathbb{Z}_{>0}$, the attention block maps \mathcal{K}_0^t into \mathcal{K}_1^t .

⁴⁴Here $(\mathbb{R}^D)^+$ is the set of nonempty tuples of vectors in \mathbb{R}^D . This is formally the disjoint union (coproduct) $\coprod_{t \in \mathbb{Z}_{>0}} \mathbb{R}^{t \times D}$.

Note that we *cannot* use Theorem 5.2.1 here because the attention block defines a different function on $\mathbb{R}^{t \times D} \rightarrow \mathbb{R}^{t \times D}$ for each t , and Theorem 5.2.1 only implies that there exists a separate *length-dependent* output compact set $\mathcal{K}_t \subset \mathbb{R}^{t \times D}$ for each t , which is different from this lemma’s statement.

The feedforward function is a continuous function on $\mathbb{R}^D \rightarrow \mathbb{R}^D$, and therefore, by Theorem 5.2.1, maps its input compact set \mathcal{K}_1 to an output compact set, which we call \mathcal{K}_2 .

Finally, residual connections and layer norms are also continuous functions acting on each of the input vectors, and hence by the same reasoning would also preserve compactness.

Now we can use induction and show that there exist compact sets $\mathcal{K}_3, \mathcal{K}_4, \dots, \mathcal{K}_{2L-1}, \mathcal{K}_{2L}$ where \mathcal{K}_{2L} contains the output set of the final layer. Set $\mathcal{K}' = \mathcal{K}_{2L}$ and we have proven the statement. ■

Now recall that a transformer language model with the softmax projection function (Definition 5.2.3) defines the conditional probabilities using the softmax transformation

$$p_{\text{SM}}(\bar{y}_t \mid \mathbf{y}_{<t}) = \frac{\exp(\mathbf{e}(\bar{y}_t)^\top \mathbf{h}_t)}{\sum_{\bar{y}' \in \bar{\Sigma}} \exp(\mathbf{e}(\bar{y}')^\top \mathbf{h}_t)} \quad (5.152)$$

where $\mathbf{e}(\bar{y}) \in \mathbb{R}^D$ is the output symbol embedding of $\bar{y} \in \bar{\Sigma}$ and \mathbf{h}_t is defined from the input embeddings of $\mathbf{y}_{<t}$ via Eq. (5.149). Using Lemma 5.2.1, together with the finiteness of the vocabulary Σ and the continuity of the softmax transformation (5.152), readily yields our main result on transformer language models.

Theorem 5.2.2. *The representation-based locally normalized language model (cf. Definition 5.2.3) defined by any (fixed-depth) transformer with soft attention is tight.*

Proof. Given the Transformer, there exists a fixed compact set \mathcal{K} that will contain all inputs $\mathbf{v}_i \in \mathbb{R}^D$ to the first layer. This is true because each \mathbf{v}_i is the sum of a word embedding, which falls in a finite set since $\bar{\Sigma}$ is finite, and a position embedding, which lies in the compact set $[-1, 1]^D$. Hence, by Lemma 5.2.1, there exists a fixed compact set \mathcal{K}' that contains all output embedding vectors (regardless of how long the sequence is).

The final output probability is given by a multiplication with the word embedding matrix followed by the softmax function as in Eq. (5.152). This process amounts to composing two continuous functions. In particular, we can extract the EOS probability as a continuous \mathbb{R} -valued function $g^{\text{EOS}} : \mathcal{K}' \rightarrow (0, 1)$ (neither 0 nor 1 is in the range of the softmax function). By continuity of g^{EOS} and Theorem 5.2.1, $\mathcal{K}'' \stackrel{\text{def}}{=} g^{\text{EOS}}(\mathcal{K}') \subseteq (0, 1)$ is compact. Since \mathcal{K}'' is compact, and hence closed, $\inf \mathcal{K}'' \in \mathcal{K}''$. Thus $\inf \mathcal{K}'' \in (0, 1)$ and in particular $\inf \mathcal{K}'' > 0$. Therefore, taking $\epsilon = \inf \mathcal{K}''$, we have shown that the EOS probability of a Transformer is bounded below by some $\epsilon > 0$ (regardless of the length of the sequence). Hence, by Proposition 2.5.6, any transformer-based sequence model is tight and thus defines a language model. ■

5.3 Computational Expressiveness of Transformers

So far, we have introduced our formal definition of the transformer architecture and examined its tightness. We now move on to the computational power of the architecture. This section mirrors §5.1.7 and examines the expressivity of the transformer language model as defined in Definition 5.2.3.

Transformers are a much more recent architecture than recurrent neural language models, and our theoretical understanding of them is thus much more limited. However, over the last few years, a series of results showing various properties of the transformer model have been established. At first glance, one might find a number of contradictions among them: One of the first results shows that transformers are not even able to recognize the very simple FIRST language recognized by the (unweighted) finite-state automaton shown in Fig. 5.22 nor the Dyck language. On the other hand, there is work showing that transformers *can* recognize the majority language (determining whether a string contains more symbols a or b) and can even *count*: Both of these languages are instances of non-regular languages. Moreover, a fundamental result by Pérez et al. (2021) even shows that transformers are *Turing complete*. Upon closer inspection, the results can be explained by the different *theoretical abstractions* of the original transformer model that the different works make and even different notions of equivalence. Even very subtle differences in the model can lead to substantial differences in the expressivity of the model, as we will see below. In this section, we present some original results which show that transformers can, with infinite precision, in fact, reach up to the top of the hierarchy of formal languages that we consider in these notes: They are Turing complete. We

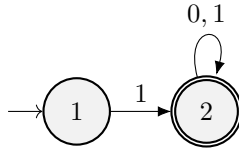


Figure 5.22: An FSA recognizing the language $\text{FIRST} = \{y \in \Sigma^* \mid \Sigma = \{0, 1\}, y_1 = 1\}$.

also comment on the differences in our approach to the work so far (the main difference and novelty is that we embed the analysis into our language modeling framework) and try to unify it. We will show that infinite-precision transformers can simulate the recognizers across the entire hierarchy: (weighted) finite-state automata, pushdown automata, and Turing machines. Luckily, apart from the first construction, the proofs and constructions in our ascent up the hierarchy will be based on a unified approach in which we build on Pérez et al. (2021) and sequentially add components to be able to recognize more and more complex languages.

Besides being more novel and thus less researched, transformers are also less intuitive to think about as sequential machines transitioning between states as with finite-state or pushdown automata. All classical computational models we introduced (finite-state automata, pushdown automata, and Turing machines) rely on some notion of an internal state which is sequentially updated, where the next state is determined based on the current configuration. We also said in §5.1.6 that this sequentiality is the Achilles’ heel of the ability to parallelize and thus speed up the computations in a language model. One of the main motivations for defining the transformer model is to avoid these sequential dependencies and to make sure the contextual representations of the individual symbols can be computed independently. However, the lack of sequentiality in transformers makes it more difficult to compare to classical and well-understood models of computation—they simply do not define any notion of a configuration that would be passed over by reading a symbol at a time, and relating the configurations at different time points to the configuration of some classical model of computation was the main idea of most of the analyses in §5.1.7. This will not be possible with transformers, and we will have to be more clever about it to draw parallels to better-understood formalisms. What is more, it seems like their parallelizable nature is one of the reasons for the lower (or, at least, ambiguous) computational power *under some formalisms*, as covered in Merrill et al. (2022a); Merrill and Sabharwal (2023).

Infinite-precision vs. finite-precision.

A word on model equivalence. As mentioned above, the nature of the transformer architecture does not lend itself well to a straightforward comparison to classical models of computation. To make the connection, we will have to be somewhat clever about the analysis. As we will see shortly, we will mainly deal with this in two ways: (i) By foregoing any notion of a state of a machine in case of n -gram language models⁴⁵ and (ii) by *embedding* the state of a computational model *into the alphabet itself*—the model will then use the augmented output alphabet to keep track of its state in the string itself without relying on any notion of its own internal state which would have to be updated sequentially.⁴⁶ How can this help us? As will become clear in our analysis of the Turing completeness of a transformer model, the model can use the *generated string* as a sort of a *sequential memory structure*. Because the transformer model can look back at the entirety of the string when computing $\text{enc}_{\mathcal{T}}(\mathbf{y}_{\leq t})$ (where $\mathbf{y}_{\leq t}$ is the augmented string generated so far), it is able to “read off” its internal state from the string. Importantly, the generated string will still contain the information about the generated string, besides including the state of the computational model. As the transformer will then compute the new embeddings $\text{enc}_{\mathcal{T}}(\mathbf{y}_{\leq t})$, it will be able to account for the state it should be in. This is illustrated in Fig. 5.23.

⁴⁵Reminder that n -gram language models are in fact *subregular* (cf. §4.1.5) and we make use of that in our analysis. Because their recognition relies purely on *local* patterns in the strings, and a transformer model has the ability to consider large enough substrings, we will see that we can model an n -gram language model without keeping any notion of a state in a transformer

⁴⁶Recall that, as discussed in §5.1.6, generation is inherently sequential. One can thus imagine augmenting the alphabet as a sort of exploitation of this sequential nature.

Original tape of the Turing machine	<table><tr><td>y_1</td><td>y_2</td><td>y_3</td><td>y_4</td><td>\dots</td><td>y_t</td><td>\dots</td></tr></table>	y_1	y_2	y_3	y_4	\dots	y_t	\dots
y_1	y_2	y_3	y_4	\dots	y_t	\dots		
Augmented tape of the stateless model	<table><tr><td>y_1q_1</td><td>y_2q_2</td><td>y_3q_3</td><td>y_4q_4</td><td>\dots</td><td>y_tq_t</td><td>\dots</td></tr></table>	y_1q_1	y_2q_2	y_3q_3	y_4q_4	\dots	y_tq_t	\dots
y_1q_1	y_2q_2	y_3q_3	y_4q_4	\dots	y_tq_t	\dots		

Figure 5.23: An abstract illustration of how a model can keep track of its internal state by “outputting” it into the generated string. By reading the augmented symbol generated at time t , the model can then determine its internal state.

While this might seem like a convenient trick to achieve Turing completeness—and in many ways, it is—it is also, in a way, cheating. This “cheating” can be described formally as the difference between model *equivalence* and *homomorphism equivalence*. When we discussed the Turing completeness of RNN LMs, we showed they can model a Turing machine by directly recognizing the same strings (for the time being, we ignored the string weights). This means that, for every Turing machine \mathcal{M} , there exists an RNN \mathcal{R} which recognizes the same language: $L(\mathcal{R}) = L(\mathcal{M})$. However, we will not be able to make statements like this in the case of transformer models. The augmented alphabet will instead bring us to a statement of the sort “For every Turing machine \mathcal{M} , there exists a transformer \mathcal{T} which recognizes the same language augmented with the state set of the Turing machine: $L_\Delta(\mathcal{T}) = L(\mathcal{M})$,” where L_Δ refers to the language of strings where each symbol is additionally augmented with the state of the Turing machine. This might seem like a small difference, but, in formal language theory, homomorphism equivalence refers to a different problem to that of normal model equivalence (Culik and Salomaa, 1978) and thus has to be considered differently. Intuitively, it additionally allows additional information to be stored in the strings (in our case, that will be the state of the Turing machine) while still considering some models to be “equivalent”. Formally, model equivalence asks the following question.

Definition 5.3.1 (Model equivalence). *Two computational models \mathcal{C}_1 and \mathcal{C}_2 are **equivalent** if*

$$L(\mathcal{C}_1) = L(\mathcal{C}_2). \quad (5.153)$$

On the other hand, homomorphic equivalence considers the following.

Definition 5.3.2. *Let \mathcal{C}_1 and \mathcal{C}_2 be two computational models. \mathcal{C}_1 is **homomorphically equivalent** to \mathcal{C}_2 if there exists a homomorphism $h: L(\mathcal{C}_1) \rightarrow L(\mathcal{C}_2)$ such that*

$$h(L(\mathcal{C}_1)) \stackrel{\text{def}}{=} \{h(\mathbf{y}) \mid \mathbf{y} \in L(\mathcal{C}_1)\} = L(\mathcal{C}_2). \quad (5.154)$$

Definition 5.3.3 (Homomorphic equivalence). *Let \mathcal{C}_1 and \mathcal{C}_2 be two computational models. \mathcal{C}_1 and \mathcal{C}_2 are **homomorphically equivalent** if \mathcal{C}_1 is homomorphically equivalent to \mathcal{C}_2 and \mathcal{C}_2 is homomorphically equivalent to \mathcal{C}_1 as per Definition 5.3.2.*

Transformers and the Inability to Recognize Simple Languages

We start our exploration of the computational power of transformer models with some negative results, which we will later “correct” by using our formalization of a transformer model and different components of the transformer architecture (for example, a different form of attention). Given their success at modeling human language which is assumed to be at least mildly context-sensitive (Huybregts et al., 1984; Shieber, 1985), it seems surprising that transformers cannot, in fact, recognize some very simple regular languages, such as PARITY or FIRST (the FSA shown in Fig. 5.22), as well as simple non-regular context-free languages such as DYCK languages:

$$\begin{aligned} \text{FIRST} &= \{\mathbf{y} \in \Sigma^* \mid \Sigma = \{0, 1\}, y_1 = 1\} \\ \text{PARITY} &= \{\mathbf{y} \in \Sigma^* \mid \Sigma = \{0, 1\}, \mathbf{y} \text{ has odd number of 1s}\} \\ \text{DYCK} &= \{\mathbf{y} \in \Sigma^* \mid \Sigma = \{(\cdot)\}, \mathbf{y} \text{ is correctly parenthesized}\} \end{aligned}$$

This has been formally shown by Hahn (2020), and experimentally verified by Chiang and Cholak (2022). Bhattamishra et al. (2020) found that transformers especially struggle to learn any languages that require counting occurrences in some way, such as the number 0s and 1s in PARITY or the number of previous open and closed parentheses in DYCK. Hahn (2020) finds that with *unique hard attention*, these languages cannot be recognized: Recognizing them by a transformer in their formulation would require the number of parameters to increase with the length of the input. Chiang and Cholak (2022) consider the setting with soft attention, where the issue is more subtle: In theory, it is possible for a transformer to recognize languages such as FIRST and PARITY, however with less and less confidence as the length increases. This is reflected by the cross-entropy of deciding language membership approaching the worst possible value of 1 bit per symbol. The reason behind this is quite intuitive: The membership of any of the languages defined above *changes* if a *single* symbol changes. However, by examining the information flow in a transformer, one can show that the corresponding information gets less and less weight relative to the length of the string due to the attention mechanism averaging over all positions.

Transformers Can Simulate n -gram Models

§5.3 showed that transformer models can struggle to recognize some of the simplest formal languages. While we did not discuss those results in detail, intuitively, they stem from the use of unique hard attention and the resulting inability to take into account all values whose keys maximize the attention scoring function. By relaxing that restriction to *averaging* hard attention, the model becomes more expressive. To show that, we begin by looking at the very simple n -gram language models, as defined in §4.1.5. By constructing, for any n -gram model, a transformer representing it, we will show the following theorem.

Theorem 5.3.1. *Let p_{LN} be an n -gram language model. Then, there exists a transformer \mathcal{T} with $L(p_{LN}) = L(\mathcal{T})$.*

Alternatively, we could say that transformers can recognize *strictly local languages* (cf. §4.1.5).

Proof. We prove the theorem by constructing, for p_{LN} , a transformer \mathcal{T} with $L(p_{LN}) = L(\mathcal{T})$. Note that we will mostly restrict the proof to the construction of the transformer, i.e., the formal definition of its parameters. The (mostly trivial) mathematical details and derivations are left as an exercise to the reader.

Recall that, by definition, an n -gram language model considers a fixed number of previous symbols to define $p_{SM}(y_t | \mathbf{y}_{<t})$ —exactly $n - 1$ of them. The constructed transformer \mathcal{T} will capture this idea with $n - 1$ *heads*, each of them attending to *exactly one* of those positions in the previous $n - 1$ positions.⁴⁷ We can then use the symbols the heads attended to (and thus identified) to identify the current n -gram and with it define the relevant conditional distribution over the next symbol. To be able to attend to the positions of interest—the ones containing the previous $n - 1$ symbols—we have to make use of appropriate positional encodings (cf. Definition 5.2.12), which will allow the model to attend to them. The idea of the construction is abstractly illustrated in Fig. 5.24.

For hopefully a better pedagogical effect, we will present this proof from the “last” part of the construction to the “first”. We, therefore, start with the final step: Assuming we have identified the appropriate n -gram $\mathbf{y}_{t-n:t-1}\bar{y}_t$, how can we encode the conditional probability distribution $p_{SM}(\bar{y}_t | \mathbf{y}_{t-n:t-1})$? The construction we use here directly mirrors the one in Minsky’s construction (cf. Lemma 5.1.3): Knowing what the individual $p_{SM}(\bar{y}_t | \mathbf{y}_{t-n:t-1})$ for $\bar{y}_t \in \bar{\Sigma}$ are (those are, as described in §4.1.5, “hard-coded”, or specified, for each n -gram separately in a look-up table), we can simply put their logits (log probabilities; in case we are using the softmax projection function) or the probabilities directly (if we are using the sparsemax projection function) into a vector and concatenate all the constructed vectors (for all possible n -grams) into a large matrix \mathbf{E} .

If we *one-hot encode* the identified n -gram by defining

$$\text{enc}_{\mathcal{T}}(\mathbf{y}_{<t}) \stackrel{\text{def}}{=} \llbracket \mathbf{y}_{t-n:t-1} \rrbracket \quad (5.155)$$

we can then, using the formulation of the transformer sequence model from Definition 5.2.3, use the one-hot encoded n -gram to lookup the appropriate column containing the conditional probabilities given the identified

⁴⁷Note that, given an n -gram model, the number n is fixed. This means that, for a given n -gram we can always fix the number of heads and therefore construct such a transformer.

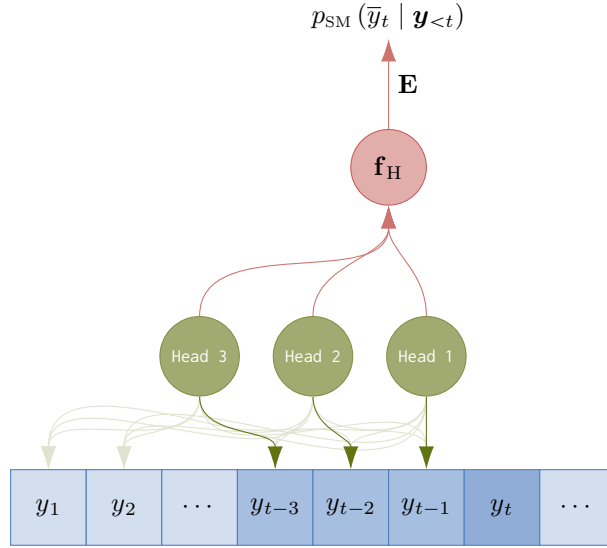


Figure 5.24: An abstract depiction of how a transformer can simulate an n -gram model using $n - 1$ heads (here, $n = 4$). The stronger arrows from the heads to the symbols in the string show where the heads concentrate their attention. The lighter green arrow is meant to represent that the heads still can consider the entire history of the input so far but are then *configured* such that they only look at the appropriate position.

n -gram for all possible $\bar{y}_t \in \bar{\Sigma}$.⁴⁸ The formal proof of correctness given that we have identified the correct n -gram is therefore analogous to the final part of the Minsky construction.

We now consider the preceding step of the simulation: How can we identify the complete n -gram given that the $n - 1$ heads of the transformer identified the symbols in the positions they attended to? This, it turns out, is a simple instance of the “AND” problem investigated in Fact 5.1.6: After concatenating the values of the $n - 1$ heads into a common vector \mathbf{v} (each of which is a $|\bar{\Sigma}|$ -dimensional vector), this vector of size $|\bar{\Sigma}|(n - 1)$ will contain the **multi-hot** representation of the n -gram of interest. Let $\bar{y}_1, \dots, \bar{y}_{n-1}$ be the symbols represented by \mathbf{v} . This means that \mathbf{v} is of the form

$$\mathbf{v} = \begin{pmatrix} \llbracket \bar{y}_1 \rrbracket \\ \vdots \\ \llbracket \bar{y}_{n-1} \rrbracket \end{pmatrix} \quad (5.156)$$

and $\mathbf{v}_{k|\bar{\Sigma}|+j} = 1$ if and only if $\bar{m}(\bar{y}_k) = j$ for an ordering \bar{m} of $\bar{\Sigma}$ determining the one-hot representations of the individual symbols. We would then like to transform this vector into a vector $\mathbf{u} \in \mathbb{R}^{|\bar{\Sigma}|^{n-1}}$ such that

$$\mathbf{u}_i = 1 \text{ if and only if } i = s(\bar{y}_1, \dots, \bar{y}_{n-1}) \quad (5.157)$$

for some ordering s of $\underbrace{\bar{\Sigma} \times \dots \times \bar{\Sigma}}_{n-1 \text{ times}}$. This can be equivalently written as

$$\mathbf{u}_i = 1 \text{ if and only if } \mathbf{v}_{k|\bar{\Sigma}|+\bar{m}(\bar{y}_k)} = 1 \text{ for all } k = 1, \dots, n - 1 \quad (5.158)$$

where $i = s(\bar{y}_1, \dots, \bar{y}_{n-1})$. Clearly, this is the same problem as described in Fact 5.1.6 and can therefore be solved by a linear transformation followed by the application of the thresholded sigmoid nonlinearity, which will together form the transformation \mathbf{f}_H combining the information obtained from all the heads of the transformer model. Note that, to make this more similar to the practice of how transformers are actually implemented, we could also use the ReLU activation function instead of the saturated sigmoid.

This brings us to the final part of the proof, which considers the first part of determining the conditional probability of the n -gram model by the transformer: Identifying the symbols at the previous $n - 1$ positions

⁴⁸Note that we are again working over the set of *extended* reals $\bar{\mathbb{R}} = \mathbb{R} \cup \{-\infty, \infty\}$ in case of the softmax activation function.

by the $n - 1$ heads of the transformer. To show how this can be done, let us consider and define the “degrees of freedom” we have left when specifying a transformer model in our framework.

- The symbol representations \mathbf{r} . We will use simple one-hot encodings of the tokens: $\mathbf{r}(\bar{y}) \stackrel{\text{def}}{=} \llbracket \bar{y} \rrbracket$.
- The positional encodings \mathbf{f}_{pos} . We will use the following simple positional encoding: $\mathbf{f}_{\text{pos}}(t) = \begin{pmatrix} t \\ 1 \end{pmatrix}$. The utility of the constant 1 will be made clear shortly. We will combine positional encodings with symbol representations by *concatenating* them into a vector of size $|\bar{\Sigma}| + 2$.
- The number of transformer layers. We will use a single transformer layer.
- The number of heads H : As we mentioned, we will use $H = n - 1$ heads to attend to the previous $n - 1$ symbols.
- The form of the attention scoring function f . While not the most typical, we will use the following scoring function:

$$f(\mathbf{q}, \mathbf{k}) \stackrel{\text{def}}{=} -|\langle \mathbf{q}, \mathbf{k} \rangle|. \quad (5.159)$$

It will, together with the positional encodings, allow us to easily single out the positions in the string that we care about.

- The form of attention. We will use hard attention (in this case, it can be either unique or averaging).
- The parameters of each of the attention heads, that is the transformations Q , K , and V . Each of those will take the form of a linear transformation of the symbol embedding. We describe them and their roles in more detail below.

As mentioned above, the input symbol \bar{y}_t is presented to the transformer model together with its positional encoding in the form

$$\mathbf{r}(\bar{y}_t) = \begin{pmatrix} \llbracket \bar{y}_t \rrbracket \\ t \\ 1 \end{pmatrix} \in \mathbb{R}^{|\bar{\Sigma}|+2}. \quad (5.160)$$

The parameters of all the heads are defined in the same way, with the only difference being a simple parameter that depends on the “index” of the head we are considering, h . Therefore, in the following, we describe the construction of a single head h . At any time step t (i.e., when modeling the conditional distribution $p_{\text{SM}}(\bar{y}_t \mid \mathbf{y}_{<t})$), the head h will attend to or be “responsible for” recognizing the symbol at position $t - h$, y_{t-h} . This can be seen in Fig. 5.24, where, for example, Head 3 is responsible for the position $t - 3$, which is denoted by the stronger arrow to that position. All we still have to do is describe the individual transformations Q_h , K_h , V_h of the head h . All of them will be *linear* transformations, i.e., matrix multiplication:

$$\mathbf{q} \stackrel{\text{def}}{=} Q(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{Q}_h \mathbf{x} \quad (5.161)$$

$$\mathbf{k} \stackrel{\text{def}}{=} K(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{K}_h \mathbf{x} \quad (5.162)$$

$$\mathbf{v} \stackrel{\text{def}}{=} V(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{V}_h \mathbf{x} \quad (5.163)$$

We now define the matrices \mathbf{Q}_h , \mathbf{K}_h , and \mathbf{V}_h , specifically in the *first* (in this case, the only) layer of a transformer language model. Importantly, since we are talking about only the first layer, we can simply consider as inputs to the layer the original static symbol representations together with their position encodings rather than any contextual representations. First, let us consider again what roles the matrices play in computing $\text{enc}_{\mathcal{T}}(\mathbf{y}_{<t})$. In the context of language modeling, the matrix \mathbf{Q}_h takes in the representation of the “latest” generated symbol y_{t-1} and produces from it the query vector of y_{t-1} . It is, therefore, only applied *once* per generation step—only for symbol y_{t-1} . The matrices \mathbf{K}_h and \mathbf{V}_h , on the other hand, transform *all* non-masked input symbols to the key and value vectors. That is, they take the representations of the input symbols and their positions encodings $\text{enc}_{\mathcal{T}}(y_j)$ for every $j = 1, \dots, t - 1$ and transform them into the key

and value vectors. The keys will then be compared with the query constructed for y_{t-1} with the \mathbf{Q}_h matrix, while the constructed values will be used to compute the new hidden state \mathbf{h}_t .⁴⁹

So, what kind of query, key, and value vectors do we want? As mentioned, the head h will be responsible for identifying the symbol at position $t - h$. Therefore, we want it to put all its attention to this position. In other words, given the query \mathbf{q}_{t-1} , we want the attention function in Eq. (5.159) to be maximized by the key of the symbol at position $t - h$. Notice that, therefore, the key does not have to depend on the identity of the symbol at position $t - h$ —only the position information matters. Let us then consider the following query and key transformations for head h :

$$Q: \begin{pmatrix} \llbracket \bar{y}_t \rrbracket \\ t \\ 1 \end{pmatrix} \mapsto \begin{pmatrix} t - h \\ 1 \end{pmatrix} \quad (5.164)$$

$$K: \begin{pmatrix} \llbracket \bar{y}_j \rrbracket \\ j \\ 1 \end{pmatrix} \mapsto \begin{pmatrix} -1 \\ j \end{pmatrix}. \quad (5.165)$$

Given such a query and such keys, the attention scoring function computes

$$f(\mathbf{q}_t, \mathbf{k}_j) = -\left| \left\langle \begin{pmatrix} t - h \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ j \end{pmatrix} \right\rangle \right| = -|t - h - j|, \quad (5.166)$$

which is maximized exactly when $j = t - h$, that is, at the position that we want the head h to attend to! This means that the hard attention we use will put all its probability mass to exactly the position we intended it to. Intuitively, both transformations keep only the positional information. The query transformation “injects” the knowledge of which position should maximize the attention score, while the key transformation (which is, again, applied to all the non-masked positions) simply “exposes” the positional information about the symbol. The alternating constant 1 (or -1) and the index of the position ensure that the inner product simply computes the *difference* between the position of the symbol and the position of interest—we will use this trick multiple times in later constructions as well. It is easy to see that the two transformations are indeed linear.

This leaves us with the question of how to use this position of the symbol of interest ($t - h$) to extract the one-hot encoding of the symbol at that position. Luckily, due to the information contained in the symbol representations $\mathbf{r}(y_j)$, this is trivial. All that the transformation V has to do is the following:

$$V: \begin{pmatrix} \llbracket \bar{y}_j \rrbracket \\ j \\ 1 \end{pmatrix} \mapsto \llbracket \bar{y}_j \rrbracket. \quad (5.167)$$

With this, the identity of the symbol is carried forward through the attention mechanism. Again, is easy to see that this is a linear transformation of the symbol representation. Notice that the only head-dependent transformation is the query transformation—it depends on the index of the head, determining the position of interest, meaning that every head defines a different query transformation, while the keys and values transformations are the same among all heads.

This concludes the proof. Fig. 5.25 again shows an illustration of the described model with all the defined components. ■

This proof establishes *the only* “concrete” result on the (lower bound of the) expressivity for transformers in the form of *model equivalence* (cf. Definition 5.3.1) that we know of. In the next subsections, we discuss how transformer-based language models can simulate more complex formal models. However, the simulation will not be as “direct” as the n -gram one, in the sense that we will have to work with modified alphabets which, as we noted above, results in a different notion of equivalence of models than what we have considered so far. We will thus not model the conditional probabilities $p_{\text{SM}}(y \mid \mathbf{y}_{<t})$, but rather the probabilities over

⁴⁹Importantly, in a multi-layer transformer, the queries would be constructed for *every* non-masked symbol and its representation (hidden state) would be updated. However, since the updated representations would not be used in the single layer case, we only have to compute the representation of the newest symbol in this case.

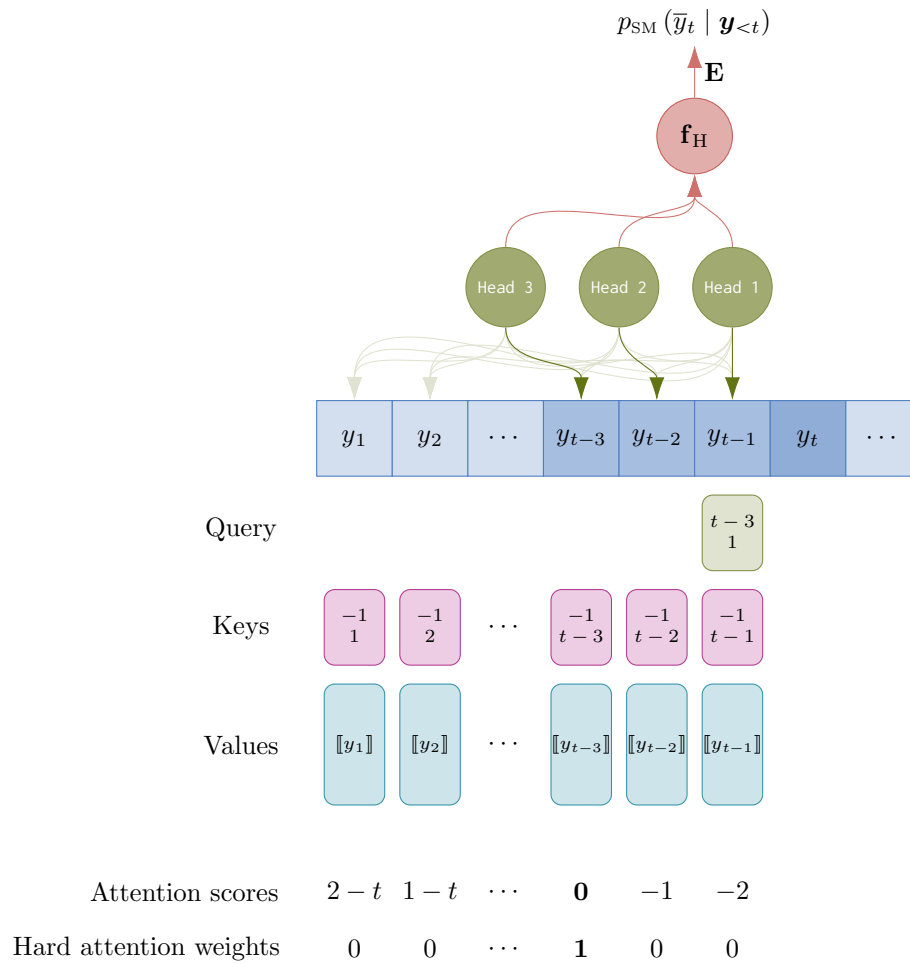


Figure 5.25: A more complete illustration of the construction described in the proof for the case of the third head, Head 3, based on Fig. 5.24. Note that, among the three heads, only the **query vector** (transformation) differs, while the key and value transformations are identical among the heads.

some more complex (but still finitely-many) objects, which will carry in them more information than just the generated symbol. As we will discuss, this will be required due to the limited abilities of transformers to execute sequential operations compared to RNNs and classical language models, as hinted at in §5.1.6.

Index

Symbols		Elman sequence model	114
bos symbol	16	embedding function	114
eos symbol	18	embedding tying	114
σ -algebra	10	empty string	13
n -gram assumption	75	energy function	16
		entropy regularizer	58
	A	event	10
accessible state	65	exploding gradient	116
activation function	112	exposure bias	52
algebra	11		
alphabet	13	F	
		finite-state automaton	60
	B	recognized language	61
bias vector	115	string acceptance	61
bigram model	75	formal language theory	13
		Frobenius normal form	74
	C		
categorical distribution	43	G	
co-accessible	65	gating functions	121
concatenation	13	generating function	92
context	18		
context-free grammar	84	H	
applicable production	84	halting problem	101
derivation	85	hidden state	107
derivation tree	85	history	18
derive	85		
language	85	I	
non-terminal	84	infinite sequence	13
parse tree	85	initial state	109
production	84	input matrix	115
start-symbol	84	input string	60
terminal	84		
context-free language	83	J	
corpus	49	Jordan sequence model	114
cross-serial dependency	106		
	D	K	
		Kleene closure	13
data leakage	54	Kleene star	13
deterministic FSA	61		
deterministic PDA	97	L	
distributed word representations	80	language	14
divergence measure	51	language model	14
dynamics map	107	context-free	91
		energy-based	16
	E	finite-state	66
early stopping	56	globally normalized	16

derivation tree weight	89	substochastic	73
induced language model	91	transition matrix	62
non-terminal allsum	90	trim	65
normalizable	90	weighted pushdown automaton	98
stringsum	89	allsum	99
weighted finite-state automaton	62	induced language model	100
allsum	65	normalizable	99
induced language model	67	run weight	98
normalizable	65	stringsum	99
state-specific allsum	65	word	13, 14
stringsum	64		

Bibliography

- Steven Abney, David McAllester, and Fernando Pereira. 1999. [Relating probabilistic grammars and automata](#). In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*, pages 542–549, College Park, Maryland, USA. Association for Computational Linguistics.
- Noga Alon, A. K. Dewdney, and Teunis J. Ott. 1991. [Efficient simulation of finite automata by neural nets](#). *J. ACM*, 38(2):495–514.
- Stéphane Aroca-Ouellette and Frank Rudzicz. 2020. [On Losses for Modern Language Models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4970–4981, Online. Association for Computational Linguistics.
- Enes Avcu, Chihiro Shibata, and Jeffrey Heinz. 2017. [Subregular complexity and deep learning](#). *ArXiv*, abs/1705.05940.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. [Layer normalization](#).
- Anton Bakhtin, Yuntian Deng, Sam Gross, Myle Ott, Marc’Aurelio Ranzato, and Arthur Szlam. 2021. [Residual energy-based models for text](#). *Journal of Machine Learning Research*, 22(40):1–41.
- Yonatan Belinkov. 2022. [Probing classifiers: Promises, shortcomings, and advances](#). *Computational Linguistics*, 48(1):207–219.
- Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. 2015. [Scheduled sampling for sequence prediction with recurrent neural networks](#). In *Advances in Neural Information Processing Systems*, volume 28.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. 2013. [Representation learning: A review and new perspectives](#). *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828. Cite arxiv:1206.5538.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, Christian Jauvin, Jauvinciro Umontreal Ca, Jaz Kandola, Thomas Hofmann, Tomaso Poggio, and John Shawe-Taylor. 2003. A neural probabilistic language model. *Journal of Machine Learning Research*, 3:1137–1155.
- Julian Besag. 1975. [Statistical analysis of non-lattice data](#). *Journal of the Royal Statistical Society. Series D (The Statistician)*, 24(3):179–195.
- Satwik Bhattamishra, Kabir Ahuja, and Navin Goyal. 2020. [On the Ability and Limitations of Transformers to Recognize Formal Languages](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7096–7116, Online. Association for Computational Linguistics.
- Patrick Billingsley. 1995. [Probability and Measure](#), 3rd edition. Wiley.
- Christopher M. Bishop. 2006. [Pattern Recognition and Machine Learning](#). Springer-Verlag, Berlin, Heidelberg.
- Mathieu Blondel, Andre Martins, and Vlad Niculae. 2019. [Learning classifiers with fenchel-young losses: Generalized entropies, margins, and algorithms](#). In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 606–615. PMLR.

- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. [Enriching word vectors with subword information](#). *Transactions of the Association for Computational Linguistics*, 5:135–146.
- L. Boltzmann. 1868. *Studien über das Gleichgewicht der lebendigen Kraft zwischen bewegten materiellen Punkten: vorgelegt in der Sitzung am 8. October 1868*. k. und k. Hof- und Staatsdr.
- T.L. Booth and R.A. Thompson. 1973. [Applying probability measures to abstract languages](#). *IEEE Transactions on Computers*, C-22(5):442–450.
- Adam L. Buchsbaum, Raffaele Giancarlo, and Jeffery R. Westbrook. 2000. [On the determinization of weighted finite automata](#). *SIAM Journal on Computing*, 30(5):1502–1531.
- Alexandra Butoi, Brian DuSell, Tim Vieira, Ryan Cotterell, and David Chiang. 2022. [Algorithms for weighted pushdown automata](#).
- Stanley F. Chen and Joshua Goodman. 1996. [An empirical study of smoothing techniques for language modeling](#). In *34th Annual Meeting of the Association for Computational Linguistics*, pages 310–318, Santa Cruz, California, USA. Association for Computational Linguistics.
- Yining Chen, Sorchia Gilroy, Andreas Maletti, Jonathan May, and Kevin Knight. 2018. [Recurrent neural networks as weighted language recognizers](#). *NAACL HLT 2018 - 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, 1:2261–2271.
- Zhiyi Chi. 1999. [Statistical properties of probabilistic context-free grammars](#). *Computational Linguistics*, 25(1):131–160.
- Zhiyi Chi and Stuart Geman. 1998. [Estimation of probabilistic context-free grammars](#). *Computational Linguistics*, 24(2):299–305.
- David Chiang and Peter Cholak. 2022. [Overcoming a theoretical limitation of self-attention](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7654–7664, Dublin, Ireland. Association for Computational Linguistics.
- Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014a. [On the properties of neural machine translation: Encoder–decoder approaches](#). In *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, pages 103–111, Doha, Qatar. Association for Computational Linguistics.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014b. [Learning phrase representations using RNN encoder–decoder for statistical machine translation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.
- N. Chomsky and M.P. Schützenberger. 1963. [The algebraic theory of context-free languages](#). In P. Braffort and D. Hirschberg, editors, *Computer Programming and Formal Systems*, volume 35 of *Studies in Logic and the Foundations of Mathematics*, pages 118–161. Elsevier.
- Noam Chomsky. 1959. [On certain formal properties of grammars](#). *Information and Control*, 2(2):137–167.
- Noam Chomsky. 1965. *Aspects of the Theory of Syntax*, 50 edition. The MIT Press.
- K. Culik and Arto Salomaa. 1978. [On the decidability of homomorphism equivalence for languages](#). *Journal of Computer and System Sciences*, 17(2):163–175.
- Grégoire Del’etang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt, Marcus Hutter, Shane Legg, and Pedro A. Ortega. 2022. [Neural networks and the chomsky hierarchy](#). *ArXiv*, abs/2207.02098.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- A. K. Dewdney. 1977. Threshold matrices and the state assignment problem for neural nets. In *Proceedings of the 8th SouthEastern Conference on Combinatorics, Graph Theory and Computing*, pages 227–245, Baton Rouge, La, USA.
- Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah A. Smith. 2020. [Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping](#). *CoRR*, abs/2002.06305.
- Li Du, Lucas Torroba Hennigen, Tiago Pimentel, Clara Meister, Jason Eisner, and Ryan Cotterell. 2022. [A measure-theoretic characterization of tight language models](#).
- John Duchi, Elad Hazan, and Yoram Singer. 2011. [Adaptive subgradient methods for online learning and stochastic optimization](#). *J. Mach. Learn. Res.*, 12(null):2121–2159.
- Rick Durrett. 2019. *Probability: Theory and Examples*, 5th edition. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge university press.
- Jason Eisner. 2016. [Inside-outside and forward-backward algorithms are just backprop \(tutorial paper\)](#). In *Proceedings of the Workshop on Structured Prediction for NLP*, pages 1–17, Austin, TX. Association for Computational Linguistics.
- Jeffrey L. Elman. 1990. [Finding structure in time](#). *Cognitive Science*, 14(2):179–211.
- Patrick C. Fischer, Albert R. Meyer, and Arnold L. Rosenberg. 1968. [Counter machines and counter languages](#). *Mathematical systems theory*, 2:265–283.
- William A. Gale and Geoffrey Sampson. 1995. [Good-turing frequency estimation without tears](#). *J. Quant. Linguistics*, 2:217–237.
- W. G. Gibbs. 1902. *Elementary Principles in Statistical Mechanics*. Charles Scribner’s Sons.
- Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2011. [Deep sparse rectifier neural networks](#). In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 315–323, Fort Lauderdale, FL, USA. PMLR.
- Glorot, Xavier and Bengio, Yoshua. 2010. [Understanding the difficulty of training deep feedforward neural networks](#). In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy. PMLR.
- Chengyue Gong, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tie-Yan Liu. 2018. [FRAGE: Frequency-Agnostic word representation](#). In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS’18, page 1341–1352, Red Hook, NY, USA. Curran Associates Inc.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.
- Andreas Griewank and Andrea Walther. 2008. *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*, 2nd edition. SIAM.
- Charles M. Grinstead and J. Laurie Snell. 1997. *Introduction to Probability*, 2nd revised edition. American Mathematical Society.
- Michael Hahn. 2020. [Theoretical limitations of self-attention in neural sequence models](#). *Transactions of the Association for Computational Linguistics*, 8:156–171.

- Paul R. Halmos and L. J. Savage. 1949. [Application of the radon-nikodym theorem to the theory of sufficient statistics](#). *The Annals of Mathematical Statistics*, 20(2):225–241.
- Yiding Hao, William Merrill, Dana Angluin, Robert Frank, Noah Amsel, Andrew Benz, and Simon Mendelsohn. 2018. [Context-free transductions with neural stacks](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 306–315, Brussels, Belgium. Association for Computational Linguistics.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. 2001. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York Inc., New York, NY, USA.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. [Delving deep into rectifiers: Surpassing human-level performance on imagenet classification](#). In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1026–1034.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. [Deep residual learning for image recognition](#). In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778.
- D.O. Hebb. 1949. *The Organization of Behavior: A Neuropsychological Theory*. A Wiley book in clinical psychology. Wiley.
- John Hewitt, Michael Hahn, Surya Ganguli, Percy Liang, and Christopher D. Manning. 2020. [RNNs can generate bounded hierarchical languages with optimal memory](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1978–2010, Online. Association for Computational Linguistics.
- John Hewitt and Christopher D. Manning. 2019. [A structural probe for finding syntax in word representations](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory](#). *Neural computation*, 9:1735–80.
- Elad Hoffer, Itay Hubara, and Daniel Soudry. 2017. [Train longer, generalize better: Closing the generalization gap in large batch training of neural networks](#). In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 1729–1739, Red Hook, NY, USA. Curran Associates Inc.
- John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. 2006. *Introduction to Automata Theory, Languages, and Computation (3rd Edition)*. Addison-Wesley Longman Publishing Co., Inc., USA.
- Roger A. Horn and Charles R. Johnson. 2012. *Matrix Analysis*, 2nd edition. Cambridge University Press.
- Ferenc Huszár. 2015. [How \(not\) to Train your Generative Model: Scheduled Sampling, Likelihood, Adversary?](#) *CoRR*, abs/1511.05101.
- Riny Huybregts, Germen de Haan, Mieke Trommelen, and Wim Zonneveld. 1984. Van periferie naar kern. *Computational Linguistics*.
- Thomas Icard. 2020a. Calibrating generative models: The probabilistic chomsky-schützenberger hierarchy. *Journal of Mathematical Psychology*, 95.
- Thomas F. Icard. 2020b. [Calibrating generative models: The probabilistic chomsky-schützenberger hierarchy](#). *Journal of Mathematical Psychology*, 95:102308.
- P. Indyk. 1995. Optimal simulation of automata by neural nets. In *STACS 95*, pages 337–348, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Gerhard Jäger and James Rogers. 2012. [Formal language theory: Refining the chomsky hierarchy](#). *Philos Trans R Soc Lond B Biol Sci*, 367(1598):1956–1970.

- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. [What does BERT learn about the structure of language?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3657, Florence, Italy. Association for Computational Linguistics.
- E. T. Jaynes. 1957. [Information theory and statistical mechanics](#). *Phys. Rev.*, 106:620–630.
- Lifeng Jin, Finale Doshi-Velez, Timothy Miller, William Schuler, and Lane Schwartz. 2018. [Depth-bounding is effective: Improvements and evaluation of unsupervised PCFG induction](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2721–2731, Brussels, Belgium. Association for Computational Linguistics.
- Michael I. Jordan. 1986. [Serial order: A parallel distributed processing approach](#). *Technical report*.
- Michael I. Jordan. 1997. [Chapter 25 - serial order: A parallel distributed processing approach](#). In John W. Donahoe and Vivian Packard Dorsel, editors, *Neural-Network Models of Cognition*, volume 121 of *Advances in Psychology*, pages 471–495. North-Holland.
- Daniel Jurafsky and James H. Martin. 2009. *Speech and Language Processing (2nd Edition)*. Prentice-Hall, Inc., USA.
- Fred Karlsson. 2007. [Constraints on multiple center-embedding of clauses](#). *Journal of Linguistics*, 43(2):365–392.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *3rd International Conference on Learning Representations*.
- Samuel A. Korsky and Robert C. Berwick. 2019. [On the computational power of rnns](#).
- Matthieu Labeau and Shay B. Cohen. 2019. [Experimenting with power divergences for language modeling](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4104–4114, Hong Kong, China. Association for Computational Linguistics.
- Daniel J. Lehmann. 1977. [Algebraic structures for transitive closure](#). *Theoretical Computer Science*, 4(1):59–76.
- Chu-Cheng Lin, Aaron Jaech, Xin Li, Matthew R. Gormley, and Jason Eisner. 2021a. [Limitations of autoregressive models and their alternatives](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5147–5173, Online. Association for Computational Linguistics.
- Chu-Cheng Lin, Aaron Jaech, Xin Li, Matthew R. Gormley, and Jason Eisner. 2021b. [Limitations of autoregressive models and their alternatives](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5147–5173, Online. Association for Computational Linguistics.
- Chu-Cheng Lin and Arya D. McCarthy. 2022. [On the uncomputability of partition functions in energy-based sequence models](#). In *International Conference on Learning Representations*.
- Nelson F. Liu, Matt Gardner, Yonatan Belinkov, Matthew E. Peters, and Noah A. Smith. 2019. [Linguistic knowledge and transferability of contextual representations](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1073–1094, Minneapolis, Minnesota. Association for Computational Linguistics.
- Christopher D. Manning, Kevin Clark, John Hewitt, Urvashi Khandelwal, and Omer Levy. 2020. [Emergent linguistic structure in artificial neural networks trained by self-supervision](#). *Proceedings of the National Academy of Sciences*, 117(48):30046–30054.

- André F. T. Martins and Ramón F. Astudillo. 2016. [From softmax to sparsemax: A sparse model of attention and multi-label classification](#). In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML'16, page 1614–1623. JMLR.org.
- Clara Meister, Elizabeth Salesky, and Ryan Cotterell. 2020. [Generalized entropy regularization or: There's nothing special about label smoothing](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6870–6886, Online. Association for Computational Linguistics.
- William Merrill. 2019. [Sequential neural networks as automata](#). In *Proceedings of the Workshop on Deep Learning and Formal Languages: Building Bridges*, pages 1–13, Florence. Association for Computational Linguistics.
- William Merrill and Ashish Sabharwal. 2023. [The parallelism tradeoff: Limitations of log-precision transformers](#). *Transactions of the Association for Computational Linguistics*, 11:531–545.
- William Merrill, Ashish Sabharwal, and Noah A. Smith. 2022a. [Saturated transformers are constant-depth threshold circuits](#). *Transactions of the Association for Computational Linguistics*, 10:843–856.
- William Merrill, Alex Warstadt, and Tal Linzen. 2022b. [Entailment semantics can be extracted from an ideal language model](#). In *Proceedings of the 26th Conference on Computational Natural Language Learning (CoNLL)*, pages 176–193, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- William Merrill, Gail Weiss, Yoav Goldberg, Roy Schwartz, Noah A. Smith, and Eran Yahav. 2020. [A formal hierarchy of RNN architectures](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 443–459, Online. Association for Computational Linguistics.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. [Efficient estimation of word representations in vector space](#).
- George A. Miller and Noam Chomsky. 1963. [Finitary models of language users](#). In D. Luce, editor, *Handbook of Mathematical Psychology*, pages 2–419. John Wiley & Sons.
- Thomas Minka. 2005. [Divergence measures and message passing](#). Technical report, Microsoft.
- Marvin Lee Minsky. 1986. *Neural Nets and the brain model problem*. Ph.D. thesis, Princeton University.
- Mehryar Mohri, Fernando Pereira, and Michael Riley. 2008. [Speech Recognition with Weighted Finite-State Transducers](#), pages 559–584. Springer Berlin Heidelberg, Berlin, Heidelberg.
- James R. Munkres. 2000. *Topology*, 2nd edition. Prentice Hall, Inc.
- Hermann Ney, Ute Essen, and Reinhard Kneser. 1994. [On structuring probabilistic dependences in stochastic language modelling](#). *Computer Speech & Language*, 8(1):1–38.
- Frank Nielsen. 2018. [What is an information projection?](#) *Notices of the American Mathematical Society*, 65:1.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. [On the difficulty of training recurrent neural networks](#). In *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*, ICML'13, page III–1310–III–1318. JMLR.org.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [GloVe: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz Kaiser, and Geoffrey E. Hinton. 2017. [Regularizing neural networks by penalizing confident output distributions](#). In *Proceedings of the International Conference on Learning Representations*.
- B.T. Polyak. 1964. [Some methods of speeding up the convergence of iteration methods](#). *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17.

- Alethea Power, Yuri Burda, Harrison Edwards, Igor Babuschkin, and Vedant Misra. 2022. Grokking: Generalization beyond overfitting on small algorithmic datasets. *CoRR*, abs/2201.02177.
- Jorge Pérez, Pablo Barceló, and Javier Marinkovic. 2021. [Attention is turing-complete](#). *Journal of Machine Learning Research*, 22(75):1–35.
- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2021. [A Primer in BERTology: What We Know About How BERT Works](#). *Transactions of the Association for Computational Linguistics*, 8:842–866.
- Halsey L. Royden. 1988. *Real Analysis*, 3rd edition. Prentice-Hall.
- Thibault Sellam, Steve Yadlowsky, Ian Tenney, Jason Wei, Naomi Saphra, Alexander D’Amour, Tal Linzen, Jasmijn Bastings, Iulia Raluca Turc, Jacob Eisenstein, Dipanjan Das, and Ellie Pavlick. 2022. [The multiBERTs: BERT reproductions for robustness analysis](#). In *International Conference on Learning Representations*.
- Claude E. Shannon. 1948. [A mathematical theory of communication](#). *The Bell System Technical Journal*, 27(3):379–423.
- Stuart M. Shieber. 1985. Evidence against the context-freeness of natural language. *Linguistics and Philosophy*, 8:333–343.
- Kumar Shridhar, Alessandro Stolfo, and Mrinmaya Sachan. 2023. [Distilling reasoning capabilities into smaller language models](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7059–7073, Toronto, Canada. Association for Computational Linguistics.
- Hava T. Siegelmann and Eduardo D. Sontag. 1992. [On the computational power of neural nets](#). In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, COLT ’92, page 440–449, New York, NY, USA. Association for Computing Machinery.
- Michael Sipser. 2013. *Introduction to the Theory of Computation*, third edition. Course Technology, Boston, MA.
- Noah A. Smith and Mark Johnson. 2007. [Weighted and probabilistic context-free grammars are equally expressive](#). *Computational Linguistics*, 33(4):477–491.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. [Dropout: A simple way to prevent neural networks from overfitting](#). *Journal of Machine Learning Research*, 15(56):1929–1958.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2015. [Rethinking the inception architecture for computer vision](#). *2016 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826.
- Gábor Szárnyas. 2020. [Graphs and matrices: A translation of "Graphok és matrixok" by Dénes König \(1931\)](#).
- Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. 2020. [oLMpics-On What Language Model Pre-training Captures](#). *Transactions of the Association for Computational Linguistics*, 8:743–758.
- Terence Tao. 2011. *An Introduction to Measure Theory*. American Mathematical Society.
- Terence Tao. 2016. *Analysis II: Third Edition*. Texts and Readings in Mathematics. Springer Singapore.
- Wilson L. Taylor. 1953. ["Cloze Procedure": A new tool for measuring readability](#). *Journalism Quarterly*, 30(4):415–433.
- L. Theis, A. van den Oord, and M. Bethge. 2016. [A note on the evaluation of generative models](#). In *4th International Conference on Learning Representations*.
- A. M. Turing. 1937. [On Computable Numbers, with an Application to the Entscheidungsproblem](#). *Proceedings of the London Mathematical Society*, s2-42(1):230–265.

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems*, volume 30.
- Gail Weiss, Yoav Goldberg, and Eran Yahav. 2018. [On the practical computational power of finite precision RNNs for language recognition](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 740–745, Melbourne, Australia. Association for Computational Linguistics.
- Sean Welleck, Ilya Kulikov, Jaedeok Kim, Richard Yuanzhe Pang, and Kyunghyun Cho. 2020. [Consistency of a recurrent language model with respect to incomplete decoding](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5553–5568, Online. Association for Computational Linguistics.
- George Kingsley Zipf. 1935. *The Psycho-Biology of Language*. Houghton-Mifflin, New York, NY, USA.