

Language Modeling with Reduced Densities

Tai-Danae Bradley¹ and Yiannis Vlassopoulos²

¹Sandbox@Alphabet, Mountain View, CA 94043, USA

²Tunnel, New York, NY 10021, USA

This work originates from the observation that today’s state-of-the-art statistical language models are impressive not only for their performance, but also—and quite crucially—because they are built entirely from correlations in unstructured text data. The latter observation prompts a fundamental question that lies at the heart of this paper: *What mathematical structure exists in unstructured text data?* We put forth enriched category theory as a natural answer. We show that sequences of symbols from a finite alphabet, such as those found in a corpus of text, form a category enriched over probabilities. We then address a second fundamental question: *How can this information be stored and modeled in a way that preserves the categorical structure?* We answer this by constructing a functor from our enriched category of text to a particular enriched category of reduced density operators. The latter leverages the Loewner order on positive semidefinite operators, which can further be interpreted as a toy example of entailment.

Contents

1 Introduction	1
1.1 Related Work	3
1.2 Acknowledgments	4
2 Modeling Probability Distributions with Density Operators	4
2.1 Understanding Reduced Densities	4
2.2 Why Densities for Language?	5
3 Assigning Reduced Densities to Words	7
4 Preserving a Preorder Structure	13
4.1 Language as a Preordered Set	14
5 Language as a Category Enriched Over Probabilities	15
5.1 Conclusion	17

1 Introduction

Statistical language models attempt to learn syntactic and semantic structure in language using the statistics of the language. Great progress has been made recently using neural network architectures [RNSS18, VSP+17, GGML15, DCLT19], although the problem of modeling the meaning of text—demonstrated, for instance, by a model’s ability to answer questions—is still unsolved. From a mathematical perspective, a number of questions remain unanswered: What mathematical structure captures the meaning of expressions in a natural language? How much of this structure can be sufficiently detected with corpora of text? Is there a way to naturally mine abstract concepts

Tai-Danae Bradley: tai.danae@math3ma.com

Yiannis Vlassopoulos: yiannis@tunnel.tech

and their interrelations? How do logic and propositional entailment arise? Even so, today’s state of the art statistical language models are quite impressive for being built only from correlations in unstructured text data. This observation prompts a fundamental question that lies at the heart of this paper: *What mathematical structure exists in unstructured text data?* We propose that enriched category theory provides a natural home for the answer. In particular we show that sequences of symbols from a finite alphabet, such as those found in a corpus of text, form a category enriched over probabilities. Category theory thus gives a principled means of organizing “what goes with what” in a corpus of text along with the statistics of the resulting expressions—precisely the information used as input to today’s statistical language models. Equipped with this understanding, we then turn to another fundamental question: *How can this information be stored and modeled in a way that preserves the categorical structure?* In other words, what is a representation of this mathematical structure? We propose the answer lies in a functor from our enriched category of text to a particular enriched category of linear operators. Unwinding the details is the primary goal of this paper.

We have been led to these tools from a number of independent, yet complementary, viewpoints on mathematical structure in natural language. On meaning, we take inspiration from the Yoneda lemma in category theory, which informally states that a mathematical object is completely determined by the network of relationships that object has with other objects in its environment. In this way, two objects are isomorphic if and only if they share the same network of relationships. Putting this in the context of language, we view the meaning of a word as captured by the network of ways that word fits into other expressions in the language. In the words of linguist John Firth, “You shall know a word by the company it keeps” [Fir57]. Distinguishing the meanings of words thus amounts to distinguishing the environments in which they occur *in addition to*, as we put forth in this paper, the statistics of these occurrences. Towards marrying these ideas, another viewpoint we take is that language exhibits both algebraic (or compositional) and statistical structure. Intuitively, language can be viewed as an algebra whose elements are expressions in the language and where the product of two expressions is their concatenation if the result is meaningful and is zero otherwise. Computing a representation of this algebra using statistical data in real-world text leads one directly to tensor networks [PTV17], which are factorizations of high-dimensional tensors often used for modeling states of complex quantum many body systems [Orú19, BC17, Bia19, Pen71, Sch11], a point revisited below.

What’s more, our algebraic viewpoint is not incompatible with our category theoretical one. For instance, the “network of ways a word fits into other expressions” may be identified with the two-sided ideal of that word. If the algebra were commutative, then we could think of language as a coordinate algebra on its spectrum whose points are the (prime) ideals of the algebra, and language would be considered as the coordinate algebra on the space of meanings (and so a translation would be a change of coordinates). A non-commutative version of this leads one to considering a category of modules, though this is something to be studied in the future. Even so, algebraic structure alone is not sufficient to understand mathematical structure in language. Statistical features also play a vital role. Indeed, language exhibits long-distance correlations decaying with a power-law, and small-scale perturbations can propagate to all scales. For instance, changing a single word in an expression can change the entire meaning of the text: *I’m going to the post office* and *I’m going postal* have drastically different meanings. Such features are also characteristic of quantum critical systems [LT17], which are efficiently modeled by certain tensor networks. Inspiration also comes from a linguistic perspective, as Chomsky’s linguistic theory of generative grammars leads to tensor networks as soon as one tries to make them probabilistic [GO19, DeG19].

The primary contributions of this work are theoretical, but let us emphasize a salient point about practical implementations. As alluded to above, the statistics in language observed in corpora of text resembles the same statistics observed in one-dimensional quantum critical systems, and the ground states of the latter are known to be well approximated by low rank tensor factorizations—see [LT17, and references therein] as well as [KMH+20, GO19, EV11, PTV17, PV17]. With this observation in mind, we therefore work under the premise that the linear operators in our framework can be approximated efficiently by tensor networks. Indeed, the ability of tensor network methods and algorithms to efficiently handle and store data in ultra large-dimensional spaces is well-understood—reviewing it here is beyond the scope of this paper, but see [RL19, Orú19, Orú14, Ose11] and references therein. Further note that such techniques have

found an increasing success in machine learning in recent years, including anomaly detection, image classification, audio classification, language modeling, and more [WRVL20, SS16, ST19, MRT21, RS20, RL19, MVRL20, GPC20, GJLP18, CSS16].

With this motivation in hand, the paper is organized as follows. Our model begins by viewing language as sequences from some finite vocabulary and by viewing the statistics of language as modeled by a probability distribution on this set. Section 2 recalls a particular passage from classical to quantum probability by modeling any probability distribution on a finite set of sequences as a particular rank 1 density operator. Section 2.2 focuses on the corresponding reduced density operators, which harness valuable statistical information about the original probability distribution. We review this information in the context of language. Section 3 builds on this framework to give the main construction, namely the assignment to any word or phrase s a particular reduced density operator ρ_s . As seen in Corollary 3.1, this operator has the property that it decomposes as a weighted sum of reduced density operators—one associated to each meaningful expression in the language that contains s —where the weights are conditional probabilities of containment. As a result, ρ_s captures something of the environment, or the “meaning,” of s in a highly principled way. An immediate consequence is that the passage $s \mapsto \rho_s$ also preserves a certain hierarchical structure that is exhibited by expressions in language. In particular, Section 4 shows that both expressions in language and their associated operators form preordered sets. The former will be given by subsequence containment of consecutive symbols, for instance: $red \leq red\ rose$. The latter will be given by the Loewner order, where we work with reduced densities $\hat{\rho}_s$ not normalized to unit trace. In addition to being a map of preorders, this assignment $s \mapsto \hat{\rho}_s$ also preserves statistics in a compatible way. Section 5 makes this statement precise using the language of category theory. We show that both the preorder of expressions s and the preorder of their associated operators $\hat{\rho}_s$ can be equipped with the structure of categories enriched over the unit interval. The main result in Theorem 5.1 concludes that the passage $s \mapsto \hat{\rho}_s$ amounts to an enriched functor between these enriched categories.

1.1 Related Work

Tensor network language models have been explored previously [ZSZ+18, ZZM+19, GO19], though to our knowledge these efforts do not seek to identify the mathematical structure in unstructured text data, nor do they ask for a faithful representation of such structure or work under the hypothesis that tensor networks are candidate representations of it. This foundational perspective is also absent from quantum language models such as [BT17, SNB13, LZSH16, CPD20, ZNS+18, LWM19]. Another line of work is [CSC10], which details a categorical compositional distributional (DisCo) framework for language. The authors of [BCLM19] build upon it to describe a density operator model for entailment using the Loewner order, while [PKCS15] use densities to model homonymy. More recently, density operators and neural methods come together in [ML20], while related works on modeling entailment with densities include [BSC15, SKB18] and references therein. Notably, DisCo models requires a choice of grammatical structure as input, which is not the case in our framework. Indeed, motivated by the success of today’s statistical language models trained only on unstructured text data, we propose to let statistics as a proxy for grammar. Lastly, the primary role of the Loewner order in our work is that it instantiates the existence of an enriched functor that preserves the mathematical structure present in corpora of text. This perspective is key to the work below and is absent from the approaches listed above.

The recipe in Equation (3) for passing from a probability distribution on a set of sequences to a rank 1 density operator appears in [BST20], where it is key to a certain tensor network generative model. The passage was further elaborated on in the context of algebraic and statistical mathematical structure in [Bra20], where a preliminary discussion of this paper’s framework appears in Section 3.4. We first learned of the idea to consider language as a preorder from Misha Gromov in [Gro15]. After finishing this work we noticed the same article also advocates for a “functor” from a “linguistic category” to a “small and simple category” [Gro15, p. 59]. Lastly, we occasionally use tensor network diagrams to illustrate certain constructions. The diagrams are much like category theorists’ string diagrams, and we assume some familiarity. For an introduction to these visual representations, see [Sto19, BB17, Orú14, Eve19] or [Bra20, Section 2.2.2].

1.2 Acknowledgments

The authors thank Maxim Kontsevich, Jacob Miller, and John Terilla for helpful discussions, as well as the anonymous referee for their valuable feedback.

2 Modeling Probability Distributions with Density Operators

Let S be a finite set, and let $V = \mathbb{C}^S$ denote the free complex vector space generated by S . Concretely, the elements of S define an orthonormal basis for V , and we will denote the basis vector associated to $s \in S$ using the same letter $s \in V$. If an ordering is chosen so that $S = \{s_1, \dots, s_d\}$, then by identifying each s_i with the i th standard basis vector in \mathbb{C}^d we have an isomorphism $V \cong \mathbb{C}^d$. This space has the usual inner product $\langle s_i, s_j \rangle$, which is equal to 1 if $i = j$ and is 0 otherwise. Each vector $v \in V$ defines a linear functional $v^*: V \rightarrow \mathbb{C}$ defined by $v' \mapsto \langle v, v' \rangle$. We denote the vector space of such linear functionals on V by $V^* := \text{hom}(V, \mathbb{C})$. Given a finite-dimensional space W , we may denote elements in the tensor product $V \otimes W$ with vw in lieu of $v \otimes w$. Note that each element wv^* of the tensor product $W \otimes V^*$ corresponds to a linear map $V \rightarrow W$ defined by $v' \mapsto w\langle v, v' \rangle$. In particular, let $\text{End}(V)$ denote the space of linear operators on V . Then for any unit vector $\psi \in V$, the vector $\psi\psi^* \in V \otimes V^*$ corresponds to an operator in $\text{End}(V)$ that maps ψ to itself and maps any vector orthogonal to ψ to 0. We will denote this orthogonal projection operator by Pr_ψ .

A *density operator*, or simply *density*, ρ on a Hilbert space is a unit-trace, positive semidefinite operator. We will denote the latter property by $\rho \geq 0$. Density operators are also called *quantum states* and may be thought of as the quantum analogues of classical probability distributions. Indeed, every density ρ on $V = \mathbb{C}^S$ defines a probability distribution $\pi_\rho: S \rightarrow \mathbb{R}$ on the set S by the *Born rule*, where the probability of an element $s \in S$ is defined by $\pi_\rho(s) := \langle \rho s, s \rangle$. These values are the diagonal entries of the matrix for ρ in the basis provided by S . They are nonnegative since ρ is positive semidefinite, and their sum is 1 since ρ has unit trace. Going in the other direction, any probability distribution $\pi: S \rightarrow \mathbb{R}$ gives rise to a density on V with the property that the Born distribution induced by it coincides with π . In fact there are multiple ways to define such a density. One could consider the maximal rank diagonal operator $\sum_s \pi(s) s s^*$, whose matrix representation contains the probabilities $\pi(s)$ along its diagonal and zeros elsewhere. In what follows, however, we focus on a rank 1 density operator—namely orthogonal projection onto a particular unit vector. Concretely, consider the following unit vector in V ,

$$\psi = \sum_{s \in S} \sqrt{\pi(s)} s \quad (1)$$

and let $\text{Pr}_\psi: V \rightarrow V$ denote the orthogonal projection operator onto ψ . This unit-trace operator is positive semidefinite and satisfies $\pi_{\text{Pr}_\psi}(s) = \langle \text{Pr}_\psi s, s \rangle = \sqrt{\pi(s)}^2 = \pi(s)$ as claimed. The assignment $\pi \mapsto \text{Pr}_\psi$ provides for us a key passage from classical to quantum probability whose significance is seen when the probability distribution being modeled is a joint distribution. We elaborate below.

2.1 Understanding Reduced Densities

Given finite-dimensional vector spaces A and B , there is an isomorphism $\text{End}(A \otimes B) \cong \text{End}(A) \otimes \text{End}(B)$, and the trace defines a pair of natural linear maps $\text{tr}_A := \text{tr} \otimes \text{id}_B$ and $\text{tr}_B := \text{id}_A \otimes \text{tr}$ called *partial traces* from the tensor product of the endomorphism spaces to each factor. The partial trace preserves both trace and positive semidefiniteness, and so any density operator $\rho: A \otimes B \rightarrow A \otimes B$ gives rise to *reduced density operators* $\rho_B := \text{tr}_A \rho: B \rightarrow B$ and $\rho_A := \text{tr}_B \rho: A \rightarrow A$. These operators may be thought of as the quantum analogues of marginal probability distributions. The analogy is especially clear when the original density ρ is the orthogonal projection onto a unit vector defined by a joint probability distribution. To see this, suppose $S = X \times Y$ for finite ordered sets $X = \{x_1, \dots, x_n\}$ and $Y = \{y_1, \dots, y_m\}$, and let $\pi: S \rightarrow \mathbb{R}$ be any joint probability distribution. As before, this defines the orthogonal projection operator $\text{Pr}_\psi: \mathbb{C}^X \otimes \mathbb{C}^Y \rightarrow \mathbb{C}^X \otimes \mathbb{C}^Y$

onto the unit vector

$$\psi = \sum_{i,a} \psi_{ia} x_i y_a \quad (2)$$

where $\psi_{ia} = \sqrt{\pi(x_i, y_a)}$ as in Equation (1). Explicitly,

$$\text{Pr}_\psi = \psi \psi^* = \left(\sum_{i,a} \psi_{ia} x_i y_a \right) \left(\sum_{j,b} \bar{\psi}_{jb} x_j^* y_b^* \right) = \sum_{\substack{i,a \\ j,b}} \psi_{ia} \bar{\psi}_{jb} x_i x_j^* y_a y_b^*. \quad (3)$$

An application of the partial trace yields the reduced density operator $\rho_Y : \mathbb{C}^Y \rightarrow \mathbb{C}^Y$, which has the following expression,

$$\rho_Y = \text{tr}_X \text{Pr}_\psi = \sum_{\substack{i,a \\ j,b}} \psi_{ia} \bar{\psi}_{jb} \text{tr}_X(x_i x_j^* y_a y_b^*) = \sum_{\substack{i,a \\ j,b}} \psi_{ia} \bar{\psi}_{jb} \text{tr}(x_i x_j^*) \cdot y_a y_b^* = \sum_{\substack{a,b \\ i}} \psi_{ia} \bar{\psi}_{ib} y_a y_b^* \quad (4)$$

where the last equality follows from $\text{tr}(x_i x_j^*) = \langle x_j, x_i \rangle$, which is 1 if $i = j$ and is 0 otherwise. Notice that the a th diagonal entry of ρ_Y is marginal probability $\pi_Y(y_a) := \sum_i \psi_{ia} \bar{\psi}_{ia} = \sum_i \pi(x_i, y_a)$, and so the diagonal of ρ_Y recovers the marginal probability distribution $\pi_Y : Y \rightarrow \mathbb{R}$ obtained from the joint distribution π .

$$\rho_Y = \begin{bmatrix} \pi_Y(y_1) & & & * \\ & \pi_Y(y_2) & & \\ & & \ddots & \\ * & & & \pi_Y(y_m) \end{bmatrix}$$

The ab th off-diagonal entry of this matrix is $(\rho_Y)_{ab} = \sum_i \sqrt{\pi(x_i, y_a) \pi(x_i, y_b)}$ which is generally nonzero and measures the extent to which y_a and y_b have common interactions with elements in X . This can be expressed in terms of Bhattacharyya coefficients, which measure the proximity of two probability distributions. Unwinding this, the *Bhattacharyya coefficient* for two probability distributions $p, q : S \rightarrow \mathbb{R}$ on a finite set S is defined by $B(p, q) := \sum_s \sqrt{p(s)q(s)}$. Putting this in the context of reduced densities, each suffix $y_a \in Y$ defines a conditional probability distribution $\pi_a : X \rightarrow \mathbb{R}$ by $x_i \mapsto \pi(x_i | y_a)$, and so the Bhattacharyya coefficient of two conditional distributions π_a and π_b is equal to

$$\sum_i \sqrt{\pi(x_i | y_a) \pi(x_i | y_b)} = \frac{1}{\sqrt{\pi_Y(y_a) \pi_Y(y_b)}} \sum_i \sqrt{\pi(x_i, y_a) \pi(x_i, y_b)} = \frac{(\rho_Y)_{ab}}{\sqrt{\pi_Y(y_a) \pi_Y(y_b)}}$$

and the off-diagonals of the reduced density are therefore $(\rho_Y)_{ab} = \sqrt{\pi_Y(y_a) \pi_Y(y_b)} B(\pi_a, \pi_b)$. This may also be written using the *Hellinger distance*, $H(p, q) := \sqrt{1 - B(p, q)}$. In much the same way, the reduced density $\rho_X = \text{tr}_Y \text{Pr}_\psi$ contains the marginal probability distribution $\pi_X : X \rightarrow \mathbb{R}$ on X along its diagonal and has additional nonzero off-diagonal entries. In both cases, the off-diagonals encode statistical information about subsystem interactions, and the spectral information of these reduced densities is akin to conditional probability, as it carries sufficient information to reconstruct the original state Pr_ψ . This idea is described in detail in [Bra20, chapter 3] and in [BST20], where understanding and harnessing this information is key to producing a tensor network generative model. In this paper, we explore the extent to which reduced densities obtained from classical probability distributions are useful in representing words and expressions in language.

2.2 Why Densities for Language?

To motivate the connection between reduced densities and language, this section gives a few elementary, yet illuminating, observations about these operators. We begin with some terminology. Given a pair $(x, y) \in X \times Y$, refer to x as the *prefix* of the pair and to y as the *suffix*. The first observation is that two suffixes have the same image under ρ_Y if and only if they share the same set of prefixes with the same probabilities (and similarly for prefixes and ρ_X).

Proposition 2.1. Let $\pi: X \times Y \rightarrow \mathbb{R}$ be a probability distribution and let ψ be the vector given in Equation (2). Suffixes y_c and y_d satisfy $\pi(x_i, y_c) = \pi(x_i, y_d)$ for all i if and only if they have the same image under $\rho_Y = \text{tr}_X \text{Pr}_\psi$.

Proof. If $\pi(x_i, y_c) = \pi(x_i, y_d)$ for all i , then

$$\rho_Y(y_c) = \sum_{i,a} \sqrt{\pi(x_i, y_a)\pi(x_i, y_c)} y_a = \sum_{i,a} \sqrt{\pi(x_i, y_a)\pi(x_i, y_d)} y_a = \rho_Y(y_d).$$

Conversely, if $\pi(x_i, y_c) \neq \pi(x_i, y_d)$ for some i , then $\rho_Y(y_c) \neq \rho_Y(y_d)$. \square

Reduced densities therefore classify suffixes (or prefixes in the case of ρ_X) that have the same environments and statistics within a language. For this reason, we refer to $\rho_Y(y_a)$ as the *ambiance vector* for the suffix y_a , as suffixes with the same ambient environment have the same image under ρ_Y .

Example 1. Consider the ordered sets $X = \{\text{big, tall, cold, chilly}\}$ and $Y = \{\text{mountain, winter}\}$ and suppose $T \subseteq X \times Y$ is the four-element subset

$$T = \{(\text{big, mountain}), (\text{tall, mountain}), (\text{cold, winter}), (\text{chilly, winter})\}.$$

Let $\pi: X \times Y \rightarrow \mathbb{R}$ be the probability distribution uniformly concentrated on T so that $\pi(x, y)$ is $1/4$ if $(x, y) \in T$ and is 0 otherwise. The vector in Equation (2) is a sum of all phrases in T , each weighted by the square root of its probability.

$$\psi = \frac{1}{2}(\text{big} \otimes \text{mountain} + \text{tall} \otimes \text{mountain} + \text{cold} \otimes \text{winter} + \text{chilly} \otimes \text{winter})$$

Following Equation (4), the matrix representations of the reduced density operators obtained from orthogonal projection onto ψ are given below.

$$\rho_X = \frac{1}{4} \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad \rho_Y = \frac{1}{4} \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

Observe that¹ the words *mountain* and *winter* share no common prefixes in T , and correspondingly $\rho_Y(\text{mountain}) \neq \rho_Y(\text{winter})$. Intuitively, these words appear in different contexts and have different meanings, and ρ_Y distinguishes them as such. On the other hand, the words *big* and *tall* have the same set of suffixes with identical probabilities, and $\rho_X(\text{big}) = \rho_X(\text{tall})$. Intuitively, these words have similar meanings because they appear in similar contexts, and ρ_X identifies them as such.

This example highlights yet another connection between reduced densities and language—namely that the entries of their matrix representations have simple, combinatorial interpretations when π is an empirical distribution. The diagonal entries of the matrix for ρ_Y in Example 1 are both 2 (momentarily ignoring the factor of $1/4$) because both *mountain* and *winter* appear twice in the subset T . Importantly, the off-diagonal entries of ρ_Y are zero, as *mountain* and *winter* have no common prefix in T . Similarly, the diagonals of ρ_X count the number of prefixes in T , and the off-diagonals count the number of shared suffixes that any pair of prefixes have in common. More generally, any subset $T \subseteq X \times Y$ can be thought of as a sampling from a corpus of text and defines an empirical probability distribution $\pi: X \times Y \rightarrow \mathbb{R}$ by $\pi(x, y) = \frac{1}{|T|}$ if $(x, y) \in T$ and $\pi(x, y) = 0$. The unit vector in Equation (2) is then $\psi = \frac{1}{\sqrt{|T|}} \sum_{(x,y) \in T} xy$, and the ab th off-diagonal entry of $\rho_Y = \text{tr}_X \text{Pr}_\psi$ is given by $\sum_i \psi_{ia} \bar{\psi}_{ib} = d/|T|$ where d is the number of prefixes $x_i \in X$ such that $(x_i, y_a) \in T$ and $(x_i, y_b) \in T$. A similar result holds for the reduced density ρ_X on prefixes. This combinatorial observation was used in [BST20] and later elaborated on in [Bra20, chapter 3], though the application to language was not emphasized there. Let us emphasize it now. In the context of language, reduced densities neatly package the statistical information contained in

¹Recall that elements of Y are identified with standard basis vectors in $\mathbb{C}^Y \cong \mathbb{C}^2$, so $\text{mountain} = [1 \ 0]^\top$ while $\text{winter} = [0 \ 1]^\top$. Similarly, $\text{big} = [1 \ 0 \ 0 \ 0]^\top$, and so on.

their off-diagonal entries in terms of prefix-suffix interactions. As Proposition 2.1 has shown, this contributes to an understanding of how words fit into a language.

We take these simple observations as motivation to further explore the extent to which reduced densities arising from classical distributions can model language. In Section 3 we assign to any word (or longer expression) s in language a reduced density operator $\hat{\rho}_s$ obtained from Pr_ψ , which will have the property that it contains algebraic and statistical information about the word's environment. This property, together with the Loewner order, is used in Section 4 to show that a simple concept hierarchy in language and the accompanying statistics are preserved under the passage $s \mapsto \hat{\rho}_s$. Section 5 describes how the preservation of this structure can be stated precisely in the language of category theory.

3 Assigning Reduced Densities to Words

To assign reduced density operators to words and expressions from a language, we start with a joint probability distribution as in Section 2. There, we considered a product of two sets, thought of as prefixes and suffixes. In this discussion, we'll consider a joint distribution on an N -fold product for $N \geq 2$. To this end, let X be a finite ordered set consisting of the basic building blocks of a language. We'll refer to elements of X as *words*, though they may be characters, symbols, words, etc. Let $S = X^{N-1} \times X$ denote the set of all sequences of length $N \geq 2$. We write S as a Cartesian product to obtain prefixes $(x_{i_{N-1}}, \dots, x_{i_2}, x_{i_1})$ and suffixes x_a so that each sequence $s \in S$ is a prefix-suffix pair $s = (x_{i_{N-1}} \cdots x_{i_2} x_{i_1}, x_a)$. As shown here, the concatenation $x_{i_{N-1}} \cdots x_{i_2} x_{i_1}$ will often be used in lieu of the tuple $(x_{i_{N-1}}, \dots, x_{i_2}, x_{i_1})$. Further, the indices of a prefix are labeled starting from right to left: the right-most index is i_1 and the left-most index is i_{N-1} . Consistent with Section 2, words comprising prefixes are labeled with i, j, \dots while suffixes are labeled with a, b, \dots . The set of suffixes X may be replaced by X^k for any $k \geq 1$, though we work with $k = 1$ for simplicity. Any subsequence of consecutive words in a prefix is called a *phrase*, and a word is a phrase of length one. With this setup in mind, suppose $\pi: S \rightarrow \mathbb{R}$ is any probability distribution and consider the unit vector $\psi = \sum_{s \in S} \psi_s s$ with $\psi_s = \sqrt{\pi(s)}$ for each s . Note that ψ lies in the tensor product $\mathbb{C}^S \cong V^{\otimes N-1} \otimes V$ where $V = \mathbb{C}^X$, and so since each basis vector s corresponds to a tensor product $s = x_{i_{N-1}} \cdots x_{i_1} x_a$ we may write

$$\psi = \sum_{i_{N-1}, \dots, i_1, a} \psi_{i_{N-1} \cdots i_1 a} x_{i_{N-1}} \cdots x_{i_1} x_a \quad (5)$$

where the coefficients are the square root of probabilities $\psi_{i_{N-1} \cdots i_1 a} = \sqrt{\pi(x_{i_{N-1}}, \dots, x_{i_1}, x_a)}$. Now consider the rank 1 density operator on $V^{\otimes N-1} \otimes V$ given by the orthogonal projection onto ψ ,

$$\text{Pr}_\psi = \psi \psi^* = \sum_{\substack{i_{N-1}, \dots, i_1, a \\ j_{N-1}, \dots, j_1, b}} \psi_{i_{N-1} \cdots i_1 a} \bar{\psi}_{j_{N-1} \cdots j_1 b} x_{i_{N-1}} \cdots x_{i_1} x_a x_{j_{N-1}}^* \cdots x_{j_1}^* x_b^*.$$

Similarly as done in Equation (4), tracing out the prefix subsystem gives the reduced density $\rho_V = \text{tr}_{V^{\otimes N-1}} \text{Pr}_\psi: V \rightarrow V$, which has the following explicit description.

$$\rho_V = \sum_{i_{N-1}, \dots, i_1, a, b} \psi_{i_{N-1} \cdots i_1 a} \bar{\psi}_{i_{N-1} \cdots i_1 b} x_a x_b^* \quad (6)$$

A slight modification of this expression gives rise to (unnormalized) reduced densities associated to phrases in X^{N-1} . Indeed, Equation (6) involves a sum over all prefixes, but suppose instead the sum is over all indices *except* those associated to a given phrase. For example, if $N = 5$ and $x_{i_2} x_{i_1}$ is a fixed phrase of length 2, consider the following operator where the indices i_1 and i_2 are fixed.

$$\hat{\rho}_{x_{i_2} x_{i_1}} := \sum_{i_4, i_3, a, b} \psi_{i_4 i_3 i_2 i_1 a} \bar{\psi}_{i_4 i_3 i_2 i_1 b} x_a x_b^*$$

This new operator may not have unit trace, though it is still positive semidefinite. We therefore view it as *the unnormalized reduced density associated to the phrase $x_{i_2} x_{i_1}$* . Informally, it is obtained by

first composing the vector $x_{i_2}x_{i_1}$ with ψ at the two sites directly adjacent to the suffix site, then forming the orthogonal projection onto this modified vector, and then tracing out the remaining prefix indices i_3 and i_4 . The tensor network diagrams in Figure 1 illustrate this relationship between ψ and ρ_V and $\hat{\rho}_{x_{i_2}x_{i_1}}$. This construction is summarized in Definition 3.1 below and an example is given in Example 2. Afterwards, we will show that renormalizing $\hat{\rho}_{x_{i_2}x_{i_1}}$ to a unit-trace operator $\rho_{x_{i_2}x_{i_1}}$ will capture the conditional probability distribution on the set of suffixes of $x_{i_2}x_{i_1}$ in a principled way.

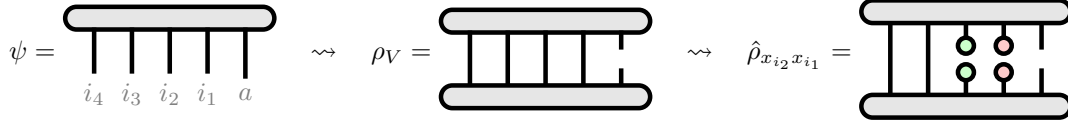


Figure 1: A tensor network diagram illustrating the construction of the reduced densities ρ_V and $\hat{\rho}_{x_{i_2}x_{i_1}}$ from the unit vector ψ .

Definition 3.1. The *unnormalized reduced density* $\hat{\rho}_{x_{i_k} \dots x_{i_1}}$ associated to a phrase $x_{i_k} \dots x_{i_1}$ of length $k \geq 1$ is the following positive semidefinite operator on $V = \mathbb{C}^X$,

$$\hat{\rho}_{x_{i_k} \dots x_{i_1}} := \sum_{\substack{i_{N-1}, \dots, i_{k+1} \\ a, b}} \psi_{i_{N-1} \dots i_1 a} \bar{\psi}_{i_{N-1} \dots i_1 b} x_a x_b^*.$$

This operator may simply be referred to as the *reduced density for* $x_{i_k} \dots x_{i_1}$, keeping in mind that it may not have unit trace.

As an immediate consequence, two words x_{i_1} and $x_{i'_1}$ map to the same operator $\hat{\rho}_{x_{i_1}} = \hat{\rho}_{x_{i'_1}}$ if they share the same statistics in the language, that is if $\psi_{i_{N-1} \dots i_1 a} = \psi_{i_{N-1} \dots i'_1 a}$ for all a and for all i_{N-1}, \dots, i_2 . The same is true for more general phrases of length $k \geq 1$. Another consequence is that the reduced density for a given phrase decomposes as a sum of other reduced densities, one associated to each expression containing that phrase. Though simple, this will result reappear a number of times.

Lemma 3.1. For any $1 \leq k \leq N - 1$ and any phrase $x_{i_k} \dots x_{i_1}$,

$$\hat{\rho}_{x_{i_k} \dots x_{i_1}} = \sum_{i_{k+1}} \hat{\rho}_{x_{i_{k+1}} x_{i_k} \dots x_{i_1}}$$

Proof. This follows directly from the definition:

$$\begin{aligned} \hat{\rho}_{x_{i_k} \dots x_{i_1}} &= \sum_{\substack{i_{N-1}, \dots, i_{k+1} \\ a, b}} \psi_{i_{N-1} \dots i_1 a} \bar{\psi}_{i_{N-1} \dots i_1 b} x_a x_b^* \\ &= \sum_{i_{k+1}} \left(\sum_{\substack{i_{N-1}, \dots, i_{k+2} \\ a, b}} \psi_{i_{N-1} \dots i_1 a} \bar{\psi}_{i_{N-1} \dots i_1 b} x_a x_b^* \right) \\ &= \sum_{i_{k+1}} \hat{\rho}_{x_{i_{k+1}} x_{i_k} \dots x_{i_1}}. \end{aligned}$$

□

As a result, the *ambience vector* associated to a word (defined below Proposition 2.1) such as $x_{i_1} = \text{dog}$ decomposes as a sum of *ambience vectors*, one for each expression ending with *dog*. The

lemma also implies² that the reduced density for any word x_{i_1} decomposes as a sum of rank 1 operators—one associated to each phrase $x_{i_{N-1}} \cdots x_{i_1}$ that ends with x_{i_1} .

$$\hat{\rho}_{x_{i_1}} = \sum_{i_2} \hat{\rho}_{x_{i_2}x_{i_1}} = \sum_{i_3, i_2} \hat{\rho}_{x_{i_3}x_{i_2}x_{i_1}} = \cdots = \sum_{i_{N-1}, \dots, i_3, i_2} \hat{\rho}_{x_{i_{N-1}} \cdots x_{i_3}x_{i_2}x_{i_1}} \quad (7)$$

To see that each operator $\hat{\rho}_{x_{i_{N-1}} \cdots x_{i_3}x_{i_2}x_{i_1}}$ has rank 1, notice that its expression in Definition 3.1 does not involve a sum over prefixes. Equivalently, no edges are contracted in its tensor diagram representation, as illustrated in Figure 2. Further observe that Definition 3.1 and Lemma 3.1 only

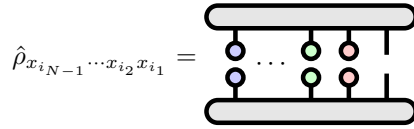


Figure 2: A rank 1 operator illustrated as a tensor network diagram with no contracted edges.

regard those phrases that occur adjacent to a suffix. One can also associate reduced densities to phrases that occur in *any* position in a sequence (for instance, see the discussion surrounding Figure 3) and find an analogous decomposition. We omit this more general discussion to streamline the presentation.

In the example below, we consider a toy corpus containing five phrases of length four. Each phrase will correspond to a sequence in a four-fold Cartesian product of *different* sets $A \times B \times C \times D$ rather than the same set $X \times X \times X \times X$, as one might use in practice. This minor adjustment will simply keep the example tidy (for instance, it allows us to consider a 2×2 matrix rather than a sparse 8×8 matrix) while still illustrating the theory.

Example 2. Begin with the following ordered sets, $A = \{\text{small, big}\}$ and $B = \{\text{black, white}\}$ and $C = \{\text{dog, cat}\}$ and $D = \{\text{barks, runs}\}$, and consider the five phrases of length $N = 4$ shown below on the left. Define the vector ψ to be the normalized sum of these five phrases, as shown on the right.

- small black dog barks
- small white dog barks
- big black dog runs
- big white cat runs
- small black cat runs

$$\psi = \frac{1}{\sqrt{5}} \begin{pmatrix} \text{small} \otimes \text{black} \otimes \text{dog} \otimes \text{barks} + \\ \text{small} \otimes \text{white} \otimes \text{dog} \otimes \text{barks} + \\ \text{big} \otimes \text{black} \otimes \text{dog} \otimes \text{runs} + \\ \text{big} \otimes \text{white} \otimes \text{cat} \otimes \text{runs} + \\ \text{small} \otimes \text{black} \otimes \text{cat} \otimes \text{runs} \end{pmatrix}$$

In fact, ψ is obtained from a probability distribution $\pi: A \times B \times C \times D \rightarrow \mathbb{R}$ as in Equation (5), where π is uniformly concentrated on the subset $T \subseteq A \times B \times C \times D$ consisting of the five phrases above. Further note that ψ lies in the tensor product $\mathbb{C}^A \otimes \mathbb{C}^B \otimes \mathbb{C}^C \otimes \mathbb{C}^D$, where each factor is isomorphic to \mathbb{C}^2 . As before, elements in $A \times B \times C$ are prefixes and elements in D are suffixes. Following Equation (6), the reduced density operator $\rho_D = \text{tr}_{A \times B \times C} \text{Pr}_\psi$ on the suffix subsystem

²In [Bra20, Section 3.4], an operator-sum decomposition of ρ_V is used to write $\hat{\rho}_{x_{i_1}}$ in terms of the linear map $V^{\otimes N-1} \rightarrow V$ associated to $\psi \in V^{\otimes N-1} \otimes V$. It is equivalent to the explicit expression given in Definition 3.1 above.

is given by

$$\begin{aligned}
\rho_D &= \sum_{(a,b,c) \in A \times B \times C} \pi(a, b, c, \text{barks}) \text{barks} \otimes \text{barks}^* \\
&+ \sum_{(a,b,c) \in A \times B \times C} \sqrt{\pi(a, b, c, \text{barks})\pi(a, b, c, \text{runs})} \text{barks} \otimes \text{runs}^* \\
&+ \sum_{(a,b,c) \in A \times B \times C} \sqrt{\pi(a, b, c, \text{runs})\pi(a, b, c, \text{barks})} \text{runs} \otimes \text{barks}^* \\
&+ \sum_{(a,b,c) \in A \times B \times C} \pi(a, b, c, \text{runs}) \text{runs} \otimes \text{runs}^* \\
&= \frac{1}{5} \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}
\end{aligned}$$

where *barks* is identified with $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ while *runs* is identified with $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$. Up to the normalizing factor $1/5$, the entries of this matrix have a simple combinatorial interpretation: *barks* appears twice in T and *runs* appears three times, and these integers are seen along the diagonal of ρ_D . The off-diagonals are both zero, which is the number of prefixes $(a, b, c) \in T$ that *barks* and *runs* have in common. To compute the reduced density $\hat{\rho}_{\text{dog}}$ associated to the word *dog* as in Definition 3.1, we fix $c = \text{dog}$ and sum over pairs $(a, b) \in A \times B$. Writing this out explicitly,

$$\begin{aligned}
\hat{\rho}_{\text{dog}} &= \sum_{(a,b) \in A \times B} \pi(a, b, \text{dog}, \text{barks}) \text{barks} \otimes \text{barks}^* \\
&+ \sum_{(a,b) \in A \times B} \sqrt{\pi(a, b, \text{dog}, \text{barks})\pi(a, b, \text{dog}, \text{runs})} \text{barks} \otimes \text{runs}^* + \\
&+ \sum_{(a,b) \in A \times B} \sqrt{\pi(a, b, \text{dog}, \text{runs})\pi(a, b, \text{dog}, \text{barks})} \text{runs} \otimes \text{barks}^* + \\
&+ \sum_{(a,b) \in A \times B} \pi(a, b, \text{dog}, \text{runs}) \text{runs} \otimes \text{runs}^* \\
&= \frac{1}{5} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}.
\end{aligned}$$

Again, the entries of this matrix can be understood combinatorially. The word *dog* appears three times in T . Of those three occurrences, it is followed by *barks* twice and by *runs* once, as seen along the diagonal of $\hat{\rho}_{\text{dog}}$. The off-diagonals are both zero, which is the number of phrases (a, b) that precede both *dog barks* and *dog runs*. Now suppose we fix $b = \text{black}$ and sum over $a \in A$ alone in the calculation above. The resulting operator is the reduced density associated to the phrase *black dog*. Further fixing $a = \text{small}$ gives the reduced density for *small black dog*.

$$\hat{\rho}_{\text{black dog}} = \frac{1}{5} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \hat{\rho}_{\text{small black dog}} = \frac{1}{5} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

Notably, operators associated to phrases containing the word *dog* are related. By Lemma 3.1 we have that $\hat{\rho}_{\text{dog}} = \hat{\rho}_{\text{black dog}} + \hat{\rho}_{\text{white dog}}$, where each summand further decomposes as $\hat{\rho}_{\text{black dog}} = \hat{\rho}_{\text{small black dog}} + \hat{\rho}_{\text{big black dog}}$ and $\hat{\rho}_{\text{white dog}} = \hat{\rho}_{\text{small white dog}}$. As a result, the reduced density for *dog* can be written as a sum of reduced densities, one for every expression containing the word *dog*. This pairs well with the intuition that the meaning of a word consists in all ways that word fits into expressions in the language.

The previous example illustrates how reduced densities obtained from a carefully chosen ψ contain relevant combinatorial and statistical information about phrases in language. Additional features of these operators are given below.

Proposition 3.1. *The trace of $\hat{\rho}_{x_{i_k} \cdots x_{i_1}}$ for any phrase $x_{i_k} \cdots x_{i_1}$ of length $k \geq 1$ is the marginal probability of that phrase.*

Proof. The trace of $\hat{\rho}_{x_{i_k} \cdots x_{i_1}}$ is the sum

$$\sum_{i_{N-1}, \dots, i_{k+1}, a} |\langle \psi, x_{i_{N-1}} \cdots x_{i_3} x_{i_2} x_{i_1} x_a \rangle|^2 = \sum_{i_{N-1}, \dots, i_{k+1}, a} \pi(x_{i_{N-1}}, \dots, x_{i_3}, x_{i_2}, x_{i_1}, x_a)$$

which is the marginal probability $\pi(x_{i_k}, \dots, x_{i_1})$. \square

As illustrated in Example 2, if π is uniformly concentrated on some subset $T \subseteq X^{N-1} \times X$, then the marginal probability $\pi(x_{i_k}, \dots, x_{i_1})$ is precisely the number of times the phrase $x_{i_k} \cdots x_{i_1}$ appears within prefixes in T , divided by the total number of sequences $|T|$. Here we omit subscripts on marginal probabilities to keep the notation clean. In any case, marginal probabilities provide a way to associate properly normalized density operators to phrases. Define the *unit-trace reduced density for a phrase* $x_{i_k} \cdots x_{i_1}$ to be the associated reduced density divided by its trace.

$$\rho_{x_{i_k} \cdots x_{i_1}} := \hat{\rho}_{x_{i_k} \cdots x_{i_1}} / \pi(x_{i_k} \cdots x_{i_1}).$$

This trace 1 operator has the property that the diagonal entries of its matrix representation in the basis given by X are the conditional probabilities $\pi(x_{i_k} \cdots x_{i_1} x_a | x_{i_k} \cdots x_{i_1})$ for each suffix $x_a \in X$. Indeed, the a th diagonal entry of $\rho_{x_{i_k} \cdots x_{i_1}}$ is $\langle \rho_{x_{i_k} \cdots x_{i_1}} x_a, x_a \rangle$ which is equal to

$$\frac{1}{\pi(x_{i_k} \cdots x_{i_1})} \sum_{i_{N-1}, \dots, i_{k+1}} |\langle \psi, x_{i_{N-1}} \cdots x_{i_1} x_a \rangle|^2 = \frac{\pi(x_{i_k} \cdots x_{i_1} x_a)}{\pi(x_{i_k} \cdots x_{i_1})} = \pi(x_a | x_{i_k} \cdots x_{i_1}).$$

In this way, the operator $\rho_{x_{i_k} \cdots x_{i_1}}$ contains the probabilities that the phrase $x_{i_k} \cdots x_{i_1}$ will be continued by a given expression.

Example 3. We resume Example 2, where the reduced density $\hat{\rho}_{\text{dog}}$ is shown below, left. The trace of this operator is $\text{tr } \hat{\rho}_{\text{dog}} = 3/5 = \pi(\text{dog}) = \sum_{a,b,d} \pi(a, b, \text{dog}, d)$, which is the number of times that *dog* appears in the toy corpus T , divided by the size of T . The unit-trace reduced density operator is below, right.

$$\hat{\rho}_{\text{dog}} = \frac{1}{5} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \quad \rho_{\text{dog}} = \frac{1}{3} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

The diagonal of ρ_{dog} is the conditional probability distribution on the set of suffixes {barks, runs} conditioned on the word *dog*. That is, $\langle \rho_{\text{dog}} \text{barks}, \text{barks} \rangle = 2/3 = \pi(\text{barks} | \text{dog})$, which is the conditional probability that a sequence $(a, b, c, d) \in T$ has $d = \text{barks}$ given that $c = \text{dog}$. Similarly, we have $\langle \rho_{\text{dog}} \text{runs}, \text{runs} \rangle = 1/3 = \pi(\text{runs} | \text{dog})$.

To motivate the next result, recall that the toy corpus in Example 2 contained two colors of dogs: *black dog* and *white dog*. As previously remarked, the meaning of the word *dog* receives a contribution from the context in which it appears—including the words *black* and *white*—together with the statistics of those appearances. The next proposition anchors this intuition on firmer ground and can be seen as an enrichment of Lemma 3.1. It states that ρ_{dog} decomposes as a weighted sum of $\rho_{\text{black dog}}$ and $\rho_{\text{white dog}}$, where the weights are conditional probabilities.

Proposition 3.2. *Let $1 \leq k \leq N-1$. The unit-trace reduced density for any phrase $x_{i_k} \cdots x_{i_1}$ can be written as a weighted sum of unit-trace reduced densities—one for each phrase of length $k+1$ ending in $x_{i_k} \cdots x_{i_1}$ —where the weights are conditional probabilities,*

$$\rho_{x_{i_k} \cdots x_{i_1}} = \sum_{i_{k+1}} \pi(x_{i_{k+1}} | x_{i_k} \cdots x_{i_1}) \rho_{x_{i_{k+1}} x_{i_k} \cdots x_{i_1}}.$$

Proof. By Lemma 3.1 we have $\hat{\rho}_{x_{i_k} \cdots x_{i_1}} = \sum_{i_{k+1}} \hat{\rho}_{x_{i_{k+1}} x_{i_k} \cdots x_{i_1}}$ and so

$$\begin{aligned} \rho_{x_{i_k} \cdots x_{i_1}} &= \frac{\hat{\rho}_{x_{i_k} \cdots x_{i_1}}}{\pi(x_{i_k} \cdots x_{i_1})} = \sum_{i_{k+1}} \frac{\pi(x_{i_{k+1}} x_{i_k} \cdots x_{i_1})}{\pi(x_{i_k} \cdots x_{i_1})} \frac{\hat{\rho}_{x_{i_{k+1}} x_{i_k} \cdots x_{i_1}}}{\pi(x_{i_{k+1}} x_{i_k} \cdots x_{i_1})} \\ &= \sum_{i_{k+1}} \pi(x_{i_{k+1}} | x_{i_k} \cdots x_{i_1}) \rho_{x_{i_{k+1}} x_{i_k} \cdots x_{i_1}}. \end{aligned}$$

\square

So Proposition 3.2 relates the unit-trace reduced density $\rho_{x_{i_k} \cdots x_{i_1}}$ to all phrases of length $k + 1$ that end with the given phrase $x_{i_k} \cdots x_{i_1}$. The proof of the following corollary gives the analogous statement for phrases of length $N - 1$.

Corollary 3.1. *The unit-trace reduced density for a word x_{i_1} decomposes as a weighted sum of rank 1 unit-trace reduced densities—one for each phrase of length $N - 1$ that ends in x_{i_1} —where the weights are conditional probabilities,*

$$\rho_{x_{i_1}} = \sum_{i_{N-1}, \dots, i_2} \pi(x_{i_{N-1}} \cdots x_{i_2} | x_{i_1}) \rho_{x_{i_{N-1}} \cdots x_{i_2} x_{i_1}}. \quad (8)$$

Proof. For any phrase $x_{i_k} \cdots x_{i_1}$ of any length $k \geq 1$, a repeated application of Proposition 3.2 shows that the unit-trace reduced density of the phrase has the following decomposition.

$$\begin{aligned} \rho_{x_{i_k} \cdots x_{i_1}} &= \sum_{i_{k+1}} \pi(x_{i_{k+1}} | x_{i_k} \cdots x_{i_1}) \rho_{x_{i_{k+1}} x_{i_k} \cdots x_{i_1}} \\ &= \sum_{i_{k+2}, i_{k+2}} \pi(x_{i_{k+2}} x_{i_{k+1}} | x_{i_k} \cdots x_{i_1}) \rho_{x_{i_{k+2}} x_{i_{k+1}} x_{i_k} \cdots x_{i_1}} \\ &= \vdots \\ &= \sum_{i_{N-1}, \dots, i_{k+1}} \pi(x_{i_{N-1}} \cdots x_{i_{k+1}} | x_{i_k} \cdots x_{i_1}) \rho_{x_{i_{N-1}} \cdots x_{i_k} \cdots x_{i_1}} \end{aligned}$$

In particular, the unit-trace reduced density associated to a given word x_{i_1} can be written as the following sum, which can be seen as an enrichment of Equation (7).

$$\rho_{x_{i_1}} = \sum_{i_{N-1}, \dots, i_2} \pi(x_{i_{N-1}} \cdots x_{i_2} | x_{i_1}) \rho_{x_{i_{N-1}} \cdots x_{i_2} x_{i_1}}$$

Recall that each $\rho_{x_{i_{N-1}} \cdots x_{i_3} x_{i_2} x_{i_1}}$ is a scalar multiple of $\hat{\rho}_{x_{i_{N-1}} \cdots x_{i_3} x_{i_2} x_{i_1}}$, and the latter has rank 1 as its explicit expression in terms of Definition 3.1 does not involve a sum over prefixes. \square

Proposition 3.2 and Corollary 3.1 model the idea that the meaning of a phrase is contained in the totality of expressions that contain it, together with the statistics of those occurrences. Notably, the decomposition in Corollary 3.1 is not unlike a “probabilistic spectral decomposition.” Indeed every self-adjoint operator, including $\rho_{x_{i_1}}$, can be written as a weighted sum of projection operators corresponding to eigenvectors. The previous corollary gives an analogous decomposition where the projections correspond not to eigenvectors but rather to vectors associated to phrases that end with the word x_{i_1} . Likewise, the weights are not eigenvalues, but are rather conditional probabilities of phrases given that their last word is x_{i_1} . Alternatively, the decomposition in Equation (8) is reminiscent of a generalized measurement in quantum mechanics—a collection of positive semidefinite operators whose sum is the identity operator. Though rather than a partition of unity, we have a partition of the observable $\rho_{x_{i_1}}$.

Example 4. Let us again resume the discussion from Example 2, where the unit-trace reduced densities associated to *black dog* and *white dog* are found to be

$$\rho_{\text{black dog}} = \frac{1}{2} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \rho_{\text{white dog}} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}.$$

Recalling the toy corpus of that example, the word *dog* appears three times. Of those three occurrences, it is preceded by *white* once and by *black* twice. Conditional probabilities are therefore $\pi(\text{white} | \text{dog}) = 1/3$ and $\pi(\text{black} | \text{dog}) = 2/3$, and indeed ρ_{dog} has the following decomposition,

$$\frac{1}{3} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} = \rho_{\text{dog}} = \pi(\text{black} | \text{dog}) \rho_{\text{black dog}} + \pi(\text{white} | \text{dog}) \rho_{\text{white dog}}$$

where the first equality was verified in Example 3. Compare this decomposition for ρ_{dog} with the unnormalized analogue $\hat{\rho}_{\text{dog}} = \hat{\rho}_{\text{white dog}} + \hat{\rho}_{\text{black dog}}$ derived in Example 2. A computation similar

to that in Example 4 shows that ρ_{dog} further decomposes into a sum of rank 1 operators,

$$\begin{aligned} \rho_{\text{dog}} &= \pi(\text{small black}|\text{dog})\rho_{\text{small black dog}} \\ &+ \pi(\text{small white}|\text{dog})\rho_{\text{small white dog}} \\ &+ \pi(\text{big black}|\text{dog})\rho_{\text{big black dog}} \\ &= \frac{1}{3} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} + \frac{1}{3} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} + \frac{1}{3} \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}. \end{aligned}$$

Compare this with the unnormalized analogue $\hat{\rho}_{\text{dog}} = \hat{\rho}_{\text{small black dog}} + \hat{\rho}_{\text{big black dog}} + \hat{\rho}_{\text{small white dog}}$ derived in Example 2.

As we'll see in Section 4, the passage from $x_{i_k} \cdots x_{i_1}$ to $\hat{\rho}_{x_{i_k} \cdots x_{i_1}}$ together with the decompositions in Lemma 3.1 and Proposition 3.2 pave the way for modeling a simple concept hierarchy in language. But first, we close this discussion with the observation that our reduced densities have an obvious left-right bias, which may be undesirable. One way to avoid this bias is simply to leave the i_{N-1} and i_{th} indices open, which gives an operator on $V \otimes V$ rather than on V alone. Figure 3 illustrates such an operator corresponding to a given phrase $x_{i_3}x_{i_2}x_{i_1}$ of length three. But whether or not the $N - 1$ st sites are kept open, let us make a final remark. Suppose for the moment that $x_{j_2}x_{j_1}$ is a phrase of length two, and suppose a corpus of text is given where both $x_{i_3}x_{i_2}x_{i_1}$ and $x_{j_2}x_{j_1}$ appear simultaneously in a longer expression. Then the decompositions of the densities $\hat{\rho}_{x_{i_3}x_{i_2}x_{i_1}}$ and $\hat{\rho}_{x_{j_2}x_{j_1}}$ in the sense of Equation (7) will contain a common summand. For example, if a corpus of text contains the expression *I walked my tiny toy poodle at the new park*, then one will find that both $\hat{\rho}_{\text{tiny toy poodle}}$ and $\hat{\rho}_{\text{new park}}$ can be written as a sum of operators, both of which will include $\hat{\rho}_{\text{I walked my tiny toy poodle at the new park}}$ as a summand. This realizes the intuition that if two different phrases are both included in a larger expression, then there is a relationship between their meanings.

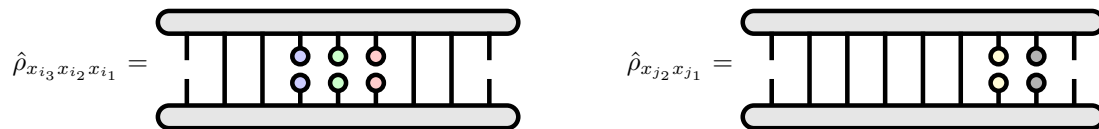


Figure 3: Associating reduced densities to phrases while leaving the first and last indices open.

4 Preserving a Preorder Structure

In this section we define a preorder on sequences and show that it is preserved under the assignment of densities to words described in the previous section. We also show it may be interpreted as a toy example of entailment. Let us begin with a definition. Given positive semidefinite operators ρ and ρ' on a fixed Hilbert space, write $\rho \geq \rho'$ if $\rho - \rho'$ is positive semidefinite. This defines a partial order on the set of such operators called the *Loewner order*. In the context of language, observe that for any phrase $x_{i_k} \cdots x_{i_1}$ and for any word $x_{i_{k+1}}$, Lemma 3.1 implies

$$\hat{\rho}_{x_{i_k} \cdots x_{i_1}} \geq \hat{\rho}_{x_{i_{k+1}} x_{i_k} \cdots x_{i_1}} \tag{9}$$

since the difference of these operators is a sum of positive semidefinite operators. For instance, in Example 2 it was shown that $\hat{\rho}_{\text{dog}} = \hat{\rho}_{\text{black dog}} + \hat{\rho}_{\text{white dog}}$ from which it follows that $\hat{\rho}_{\text{dog}} \geq \hat{\rho}_{\text{black dog}}$ and $\hat{\rho}_{\text{dog}} \geq \hat{\rho}_{\text{white dog}}$. Taking this a step further, we have

$$\hat{\rho}_{\text{dog}} \geq \hat{\rho}_{\text{black dog}} \geq \hat{\rho}_{\text{small black dog}}. \tag{10}$$

The Loewner order thus models the notion that *dog* is a more general concept than *black dog*, which is more general than *small black dog*. Just as essential, though, are the likelihood or entailment

strengths of these implications. Proposition 3.2 naturally suggests conditional probabilities as a candidate for such a measurement. Indeed, the proposition promotes Inequality (9) to the following “enriched” version,

$$\rho_{x_{i_k} \cdots x_{i_1}} \geq \pi(x_{i_{k+1}} \mid x_{i_k} \cdots x_{i_1}) \rho_{x_{i_{k+1}} x_{i_k} \cdots x_{i_1}}. \tag{11}$$

From Example 4, for instance, we find that

$$\rho_{\text{dog}} \geq \pi(\text{black} \mid \text{dog}) \rho_{\text{black dog}} \geq \pi(\text{small black} \mid \text{dog}) \rho_{\text{small black dog}}.$$

We have therefore defined a mapping from phrases to reduced densities with the property that if a phrase s' contains a phrase s of smaller length, then the corresponding operators satisfy $\rho_s \geq \pi(s' \mid s) \rho_{s'}$, or in the unnormalized case, $\hat{\rho}_s \geq \hat{\rho}_{s'}$. In short, the assignment $s \mapsto \hat{\rho}_s$ preserves a certain hierarchy that exists in the domain, namely subsequence containment. Said differently, expressions in a language form a preordered set (that is, a set equipped with a relation \leq that is reflexive and transitive), and the assignment $s \mapsto \hat{\rho}_s$ respects the preorder.

4.1 Language as a Preordered Set

The opening remarks of Section 1 shared the perspective that the meaning of a word or phrase is determined by its environment—namely, the network of ways it is contained within expressions in a language together with the statistics of those occurrences. Putting this Yoneda-lemma-like approach together with the sequential nature of language, we model the inclusion of phrases by subsequence containment. Explicitly, let X be a finite set of words, thought of as the atomic vocabulary of a language. For $N \geq 1$, let L denote the subset of $\sqcup_{k=1}^{N-1} X^k \times X$ consisting of sequences $s = (x_{i_k}, \dots, x_{i_1}, x_a)$ of all lengths $k \leq N - 1$ such that the concatenation $x_{i_k} \cdots x_{i_1} x_a$ is a meaningful expression in the language. We will refer to elements of L as *phrases*, as usual writing $x_{i_k} \cdots x_{i_1} x_a$ in lieu of $(x_{i_k}, \dots, x_{i_1}, x_a)$. If X is the set of English words, for example, then some phrases in L include *dog*, *black dog*, *tall mountain*, and *iced tea on a hot summer day*. To compare phrases $s, s' \in L$, write $s \leq s'$ if s' contains s as a subsequence. As a simple example, we write *dog* \leq *small black dog*. Next, observe that for any $s \in L$ one has $s \leq s$. Moreover, for any $s, s', s'' \in L$ if $s \leq s'$ and $s' \leq s''$, then s'' contains s' —and hence s —as a subsequence, and so $s \leq s''$. This proves the following proposition.

Proposition 4.1. *The set L equipped with the relation \leq is a preordered set.*

Some examples of comparable phrases in English include the following.

$$\begin{aligned} \text{I climbed} &\leq \text{I climbed the tall mountain} \\ \text{a hot summer} &\leq \text{iced tea on a hot summer day} \\ \text{dog} &\leq \text{black dog} \leq \text{small black dog} \end{aligned}$$

The similarity between $\text{dog} \leq \text{black dog} \leq \text{small black dog}$ and the nearly identical string of inequalities in (10) reveals an unmistakable correspondence between our preorder on language L and the Loewner order, which is a preorder on positive semidefinite operators. This correspondence is stated precisely in Proposition 4.2 below. Let us recall the setup first. Sequences in $X^N \cong X^{N-1} \times X$ are considered as prefix-suffix pairs; the initial ingredient is a probability distribution $\pi: X^{N-1} \times X \rightarrow \mathbb{R}$, which defines the unit vector $\psi \in V^{\otimes N-1} \otimes V$ in Equation (5) where $V = \mathbb{C}^X$; the vector gives rise to reduced densities of the form $\hat{\rho}_s$ or ρ_s , which correspond to (sub)sequences $s \in X^{N-1}$; and these reduced densities operate on the Hilbert space V generated by suffixes. Importantly, the assignment $s \mapsto \hat{\rho}_s$ concerns only those phrases in L of the form $s = x_{i_k} \cdots x_{i_1}$ for some $1 \leq k \leq N - 1$ that appear adjacent to a suffix, as indicated below:

$$(x_{i_{N-1}}, \dots, \overbrace{x_{i_k}, \dots, x_{i_1}}^s, x_a).$$

In what follows, we use the calligraphic font $\mathcal{L} \subseteq L$ to denote the subset of all such s .

Proposition 4.2. *Let $(\text{Pos}(V), \leq)$ denote the set of positive semidefinite operators on $V = \mathbb{C}^X$ equipped with the Loewner order, and let $\psi \in V^{\otimes N-1} \otimes V$ be the unit vector in Equation (5). The function $(\mathcal{L}, \leq) \rightarrow (\text{Pos}(V), \leq)$ defined by $s \mapsto \hat{\rho}_s$ described in Definition 3.1 is order-reversing. That is, $\hat{\rho}_s \geq \hat{\rho}_{s'}$ whenever $s \leq s'$.*

Proof. Suppose $s \leq s'$ so that $s = x_{i_k} \cdots x_{i_1}$ and $s' = x_{i_m} \cdots x_{i_{k+1}} s$ for some $1 \leq k \leq m \leq N - 1$. Lemma 3.1 implies that $\hat{\rho}_s \geq \hat{\rho}_{s'}$. \square

Observe that properties of the mapping $s \mapsto \hat{\rho}_s$ depend on the probability distribution π used to define the vector ψ . In the toy scenario of Example 2, for instance, the mapping is not “full” in the category theoretical sense since one finds that $\hat{\rho}_{\text{black cat}} \leq \hat{\rho}_{\text{dog}}$ under the Loewner order, as $\hat{\rho}_{\text{black cat}} = \frac{1}{5} \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$, and yet $\text{dog} \not\leq \text{black cat}$ in \mathcal{L} . This may not be the case, however, for different choices of π . On the other hand, $\hat{\rho}_{\text{black cat}} \leq \hat{\rho}_{\text{dog}}$ *does* imply a relationship between the sets of possible suffixes for these two phrases—see the discussion in Section 5.1. Either way, Proposition 4.2 shows that our preorder on language models a simple form of hierarchy which is preserved under the passage to linear algebra described in Section 3. Thinking back to the discussion of meaning, it is now simple to identify the context, or environment, of a word. It is simply upper closure. For instance, the set of all expressions that contain *dog* is given by³ $\uparrow \{\text{dog}\} := \{s \in L \mid \text{dog} \leq s\}$. But for simplicity, let us restrict attention to only those expressions $x_{i_k} \cdots x_{i_2} x_{i_1} \in \mathcal{L} \subseteq L$ where the last word is fixed at $x_{i_1} = \text{dog}$. In this case, the upper closure of *dog* consists of all phrases in \mathcal{L} of length at most $N - 1$ that contain *dog* as the last word. In other words, $\uparrow \{\text{dog}\}$ is equal to

$$\{x_{i_2} \text{dog} \mid x_{i_2} \in X\} \sqcup \{x_{i_3} x_{i_2} \text{dog} \mid x_{i_3}, x_{i_2} \in X\} \sqcup \cdots \sqcup \{x_{i_{N-1}} \cdots x_{i_3} x_{i_2} \text{dog} \mid x_{i_{N-1}}, \dots, x_{i_3}, x_{i_2} \in X\}.$$

As implied early on by Equation (7), the passage $\mathcal{L} \rightarrow \text{Pos}(V)$ gives rise to an analogous decomposition of words associated to these expressions:

$$\hat{\rho}_{\text{dog}} = \sum_{i_2} \hat{\rho}_{x_{i_2} \text{dog}} = \sum_{i_3, i_2} \hat{\rho}_{x_{i_3} x_{i_2} \text{dog}} = \cdots = \sum_{i_{N-1}, \dots, i_3, i_2} \hat{\rho}_{x_{i_{N-1}} \cdots x_{i_3} x_{i_2} \text{dog}}.$$

In this way, something of the “meaning” of *dog*—that is, the environment in which the word appears—is neatly packed into the single operator $\hat{\rho}_{\text{dog}}$. But as previously noted, the statistics accompanying these appearances are also essential for capturing meaning. For instance, the conditional probability of *black* given that the next word is *dog* contributes to the latter’s meaning. We therefore wish to “decorate” the preorder structure \leq with probabilities in such a way that the order-preserving map $\mathcal{L} \rightarrow \text{Pos}(V)$ retains knowledge of these conditional probabilities.

$$\text{dog} \leq \text{black dog} \quad \rightsquigarrow \quad \text{dog} \stackrel{\pi(\text{black}|\text{dog})}{\leq} \text{black dog}$$

These conditional probabilities arise naturally in the discussion on unit-trace reduced densities as in Inequality (10). For any s and s' in \mathcal{L} , the containment $s \leq s'$ implies $\rho_s \geq \pi(s' \mid s) \rho_{s'}$. But unit trace operators are not comparable under the Loewner order, and so we look for another way to incorporate probabilities with the preorder structure. We needn’t look far, however. Category theory provides a natural setting for these ideas [Rie17, Lei14, Awo10]. Indeed, every preordered set is an example of a category, and the function in Proposition 4.2 is a contravariant functor. The ability to “decorate” \leq with probabilities is the notion behind enriched category theory [Kel82], [Ell17, Appendix A]. This is made precise in Theorem 5.1 below, which generalizes Proposition 4.2 by incorporating probabilities in the desired way. We unwind this in the next section.

5 Language as a Category Enriched Over Probabilities

The full machinery of enriched category theory is not needed for our framework, and so the discussion will be kept simple. Indeed, the only categories being considered are preordered sets, and the only category we wish to enrich over is a particular symmetric monoidal preorder.

Definition 5.1. A *symmetric monoidal preorder* $(P, \leq, \cdot, 1)$ is a preorder (P, \leq) together with

- an element $1 \in P$ called the *monoidal unit*

³In the language of category theory, the upper closure of *dog* is the image of *dog* under the Yoneda embedding $L^{\text{op}} \rightarrow UL$, where UL denotes all upward closed subsets of L ordered by inclusion, and where the preorders L and UL are viewed as categories enriched over truth values.

- a function $\cdot : P \times P \rightarrow P$ called the *monoidal product*.

Moreover these data must satisfy the following properties (where we write pq for $p \cdot q$):

- for all $p, p', q, q' \in P$, if $p \leq p'$ and $q \leq q'$ then $pq \leq p'q'$,
- $1p = p1 = p$ for all $p \in P$,
- $(pq)r = p(qr)$ for all $p, q, r \in P$
- $pq = qp$ for all $p, q \in P$

The main example to have in mind is the unit interval $[0, 1] \subseteq \mathbb{R}$ equipped with the usual ordering \leq , where the monoidal product is multiplication of real numbers, and the monoidal unit is 1. In fact, $[0, 1]$ can be given the structure of a closed symmetric monoidal preorder [Eli17, Proposition 2.1.12], although we won't need closure here.

Definition 5.2. Let $(\mathcal{V}, \leq, \cdot, 1)$ be a symmetric monoidal preorder. A \mathcal{V} -enriched category \mathcal{C} , or simply \mathcal{V} -category, consists of the following data,

- a set $\text{ob}(\mathcal{C})$ of objects c, d, \dots
- an element $\mathcal{C}(c, d) \in \mathcal{V}$ for every pair of objects c and d .

Moreover, these data must satisfy the following requirements:

- $1 \leq \mathcal{C}(c, c)$ for each object c
- $\mathcal{C}(c, d) \cdot \mathcal{C}(d, e) \leq \mathcal{C}(c, e)$ for every triple of objects c, d, e .

There is also a notion of maps between \mathcal{V} -categories.

Definition 5.3. Let \mathcal{C} and \mathcal{D} be \mathcal{V} -categories. A \mathcal{V} -functor is a function $f: \text{ob}(\mathcal{C}) \rightarrow \text{ob}(\mathcal{D})$ satisfying

$$\mathcal{C}(c, c') \leq \mathcal{D}(fc, fc')$$

for all objects c, c' in \mathcal{C} .

The main result below is that both language and positive semidefinite operators form $[0, 1]$ -categories in the desired way, and moreover there is a $[0, 1]$ -functor between them. The setup is the same as before. Let X be a finite set, let $\pi: X^{N-1} \times X \rightarrow \mathbb{R}$ be any probability distribution, and let $\mathcal{L} \subseteq L \subseteq \sqcup_{k=1}^{N-1} X^k \times X$ be the subset of phrases defined in Section 4.

Proposition 5.1. *The set \mathcal{L} together with the following assignment for each $s, s' \in \mathcal{L}$ is a $[0, 1]$ -category:*

$$\mathcal{L}(s, s') = \begin{cases} \pi(s'|s) & \text{if } s \leq s' \\ 0 & \text{otherwise.} \end{cases}$$

Proof. Observe that $1 \leq \mathcal{L}(s, s)$ for each s , and it is simple to check that $\mathcal{L}(s, s')\mathcal{L}(s', s'') \leq \mathcal{L}(s, s'')$ for all s, s' and s'' . \square

With this choice of enrichment,⁴ there is a ‘‘morphism’’ between two expressions only if one is contained in the other, and moreover that morphism is labeled with the conditional probability of containment. By a similar argument, the trace on positive semidefinite operators gives rise to a $[0, 1]$ -category structure on operators assigned to phrases $s \in \mathcal{L}$. In what follows, let $\mathcal{D} \subseteq \text{Pos}(V)$ denote the image of the function $\mathcal{L} \rightarrow \text{Pos}(V)$ defined by $s \mapsto \hat{\rho}_s$.

Proposition 5.2. *The set \mathcal{D} together with the following assignment for each $\hat{\rho}_s, \hat{\rho}_{s'} \in \mathcal{D}$ is a $[0, 1]$ -category:*

$$\mathcal{D}(\hat{\rho}_s, \hat{\rho}_{s'}) = \begin{cases} \text{tr } \hat{\rho}_{s'} / \text{tr } \hat{\rho}_s & \text{if } s \leq s' \\ 0 & \text{otherwise.} \end{cases}$$

⁴The notation $\pi(s' | s)$ is used as shorthand for the conditional probability $\pi(ts|s)$ whenever $s' = ts$ for some phrase t . For example, in this section we use $\pi(\text{black dog} | \text{dog})$ to denote the conditional probability $\pi(\text{black} | \text{dog}) = \pi(\text{black dog})/\pi(\text{dog})$.

Proof. Observe that $1 \leq \mathcal{D}(\hat{\rho}_s, \hat{\rho}_s)$ for each s , and it is simple to check that $\mathcal{D}(\hat{\rho}_s, \hat{\rho}_{s'})\mathcal{D}(\hat{\rho}_{s'}, \hat{\rho}_{s''}) \leq \mathcal{D}(\hat{\rho}_s, \hat{\rho}_{s''})$ for all s, s' and s'' in \mathcal{L} . \square

Recall from Proposition 3.1 that for each $s \in \mathcal{L}$ the trace of the reduced density $\hat{\rho}_s$ is marginal probability, $\text{tr } \hat{\rho}_s = \pi(s)$. As a result, $\text{tr } \hat{\rho}_{s'} / \text{tr } \hat{\rho}_s = \pi(s') / \pi(s) = \pi(s'|s)$ whenever $s \leq s'$ and so $\mathcal{L}(s, s') \leq \mathcal{D}(\hat{\rho}_s, \hat{\rho}_{s'})$ for all $s, s' \in \mathcal{L}$. This proves the following theorem, which can be seen as an enrichment of Proposition 4.2.

Theorem 5.1. *The function $\mathcal{L} \rightarrow \mathcal{D}$ defined by $s \mapsto \hat{\rho}_s$ is a $[0, 1]$ -functor.*

We have therefore found a suitable home for modeling the passage from statistics in language to linear operators. To make this conclusion explicit, we quickly revisit the remarks given towards the end of Section 4. If s and s' are phrases in language satisfying $s \leq s'$, then the enriched categorical structure is given by conditional probability, $\mathcal{L}(s, s') = \pi(s'|s)$. Recall from Inequality (11) that this same probability arises in the relationship between the unit-trace reduced density operators associated to each phrase, $\rho_s \geq \pi(s'|s)\rho_{s'}$. Theorem 5.1 crystallizes the precise way in which these ideas connect. Indeed, if $s \leq s'$ then $\hat{\rho}_{s'} \leq \hat{\rho}_s$ which implies that $\mathcal{D}^{\text{op}}(\hat{\rho}_{s'}, \hat{\rho}_s) := \mathcal{D}(\hat{\rho}_s, \hat{\rho}_{s'}) = \text{tr } \hat{\rho}_{s'} / \text{tr } \hat{\rho}_s = \pi(s'|s)$, where the “op” denotes the opposite $[0, 1]$ -category of \mathcal{D} . This recovers the intuitive idea of decorating the network of relationships between phrases in language with the appropriate conditional probabilities. The linear algebra in Section 3 thus prescribes a principled method of assigning text to (unnormalized) reduced density operators that preserves a simple hierarchical structure in language as well as the statistics therein.

$$\begin{array}{ccc} s & & \hat{\rho}_s \\ \pi(s'|s) \downarrow & \longmapsto & \uparrow \pi(s'|s) \\ s' & & \hat{\rho}_{s'} \end{array}$$

Since the image of the functor $\mathcal{L} \rightarrow \mathcal{D}$ has more structure than its domain, we expect that operators ρ_s associated to phrases s carry additional information about the language, in addition to the simple form of entailment modeled here. The spectral information of ρ_s , for instance, may be one such source.

5.1 Conclusion

The hierarchy modeled in this work comes directly from the sequential structure of language. So while we can model the notion that *dog* is a more general concept than *small black dog*, the theory does not yet provide a way to compare phrases that aren't comparable under the preorder in \mathcal{L} . For example, one may not conclude that *mammal* abstracts the notion of *dog* since *mammal* does not contain *dog* as a subsequence. But considering all expressions that contain both *mammal* and *dog*, as in the discussion surrounding Figure 3, suggests a relationship between them. Exploring more complex hierarchies of this type is left for future work. In this direction, we make the observation that if $\hat{\rho}_{s'} \leq \hat{\rho}_s$ then for any suffix x , we have that $\langle \hat{\rho}_{s'} x, x \rangle = \pi(s'x) \leq \pi(sx) = \langle \hat{\rho}_s x, x \rangle$. In particular, if $\pi(s'x) > 0$ for some suffix x , then $\pi(sx) > 0$ as well, which is to say that any valid continuation on the right of s' is also a valid continuation of s . (In Example 2, for instance, one sees that $\hat{\rho}_{\text{black cat}} \leq \hat{\rho}_{\text{dog}}$, and that the set of suffixes of *black cat*, namely $\{\text{runs}\}$, is a subset of the set of suffixes of *dog*.) So while s and s' may not be comparable under subsequence containment, there is a clear relationship between their “right ideals,” to borrow from the algebraic perspective. Understanding this algebraic connection is left for future work. There are also additional opportunities to expand the framework using ideas from category theory. For instance, the function $x \mapsto -\log(x)$ provides a mapping $[0, 1] \rightarrow [0, \infty]$, suggesting that our framework may be reinterpreted in the theory of generalized metric spaces [Law73, Law86, Wil13]. Finally, the theory proposed in this paper fits into a larger investigation of language modeling with tensor networks. One quickly notices that the dimension of the vector spaces involved grow exponentially with the size of the vocabulary X . Realizing and manipulating ψ on a computer thus quickly becomes infeasible for real-world datasets. A similar sentiment holds for our reduced densities such as $\hat{\rho}_{x_{i_1}}$, which may operate on an ultra large-dimensional space. As discussed in Section 1, we propose these issues may be addressed by approximating ψ by a tensor network

factorization ψ_{net} , which should be chosen such that it can be computed efficiently, can faithfully model the statistics of language, and can easily generalize from a training corpus to unseen samples. We view the tensor network language models of [PTV17, PV17] as source of inspiration in this direction.

References

- [Awo10] Steve Awodey. *Category Theory*. Oxford University Press, 2010. <http://doi.org/10.1093/acprof:oso/9780198568612.001.0001>.
- [BB17] Jacob Biamonte and Ville Bergholm. Tensor networks in a nutshell. *arXiv preprint arXiv:1708.00006*, 2017.
- [BC17] Jacob C. Bridgeman and Christopher T. Chubb. Hand-waving and interpretive dance: an introductory course on tensor networks. *Journal of Physics A: Mathematical and Theoretical*, 50(22):223001, 05 2017. <http://doi.org/10.1088/1751-8121/aa6dc3>.
- [BCLM19] Dea Bankova, Bob Coecke, Martha Lewis, and Dan Marsden. Graded hyponymy for compositional distributional semantics. *Journal of Language Modelling*, 6(2):225–260, Mar. 2019. <https://doi.org/10.15398/jlm.v6i2.230>.
- [Bia19] Jacob Biamonte. Lectures on quantum tensor networks. *arXiv preprint arXiv:1912.10049*, 2019.
- [Bra20] Tai-Danae Bradley. At the interface of algebra and statistics. *arXiv preprint arXiv:2004.05631*, 2020. PhD thesis, CUNY Graduate Center.
- [BSC15] Esma Balkir, Mehrnoosh Sadrzadeh, and Bob Coecke. Distributional sentence entailment using density matrices. In *Proceedings of the First International Conference on Theoretical Topics in Computer Science*, volume 9541, pages 1–22, 2015. http://doi.org/10.1007/978-3-319-28678-5_1.
- [BST20] Tai-Danae Bradley, E. Miles Stoudenmire, and John Terilla. Modeling sequences with quantum states: A look under the hood. *Machine Learning: Science and Technology*, 1(3), 2020. <https://doi.org/10.1088/2632-2153/ab8731>.
- [BT17] Ivano Basile and Fabio Tamburini. Towards quantum language models. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1840–1849, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D17-1196>.
- [CPD20] Yiwei Chen, Yu Pan, and Daoyi Dong. Quantum language model with entanglement embedding for question answering. *arXiv preprint arXiv:2008.09943*, 2020.
- [CSC10] Bob Coecke, Mehrnoosh Sadrzadeh, and Stephen Clark. Mathematical foundations for a compositional distributional model of meaning. *Linguistic Analysis*, 36, 2010. Available at <https://arxiv.org/abs/1003.4394>.
- [CSS16] Nadav Cohen, Or Sharir, and Amnon Shashua. On the expressive power of deep learning: A tensor analysis. In Vitaly Feldman, Alexander Rakhlin, and Ohad Shamir, editors, *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pages 698–728, Columbia University, New York, New York, USA, 2016. PMLR.
- [DCLT19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 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, June 2019. Association for Computational Linguistics. <http://doi.org/10.18653/v1/N19-1423>.
- [DeG19] E. DeGiuli. Random language model. *Physical Review Letters*, 122(12), Mar 2019. <https://doi.org/10.1103/PhysRevLett.122.128301>.
- [Ell17] Jonathan Elliott. On the fuzzy concept complex. 2017. PhD thesis, University of Sheffield.
- [EV11] Glen Evenbly and Guifre Vidal. Tensor network states and geometry. *Journal of Statistical Physics*, 145(4):891–918, 2011. <https://doi.org/10.1007/s10955-011-0237-4>.
- [Eve19] Glen Evenbly. Tensors.net, 2019. Available online: <http://www.tensors.net>.

- [Fir57] John R. Firth. *A synopsis of linguistic theory 1930–55*, volume 1952–59. The Philological Society, Oxford, 1957. Reprinted in: Palmer, F. R. (ed.) (1968). *Selected Papers of J. R. Firth 1952–59*, pages 168–205. Longmans, London.
- [GGML15] Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. MADE: Masked autoencoder for distribution estimation. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 881–889, 2015.
- [GJLP18] Chu Guo, Zhanming Jie, Wei Lu, and Dario Poletti. Matrix product operators for sequence-to-sequence learning. *Physical Review E*, 98:042114, Oct 2018. <https://link.aps.org/doi/10.1103/PhysRevE.98.042114>.
- [GO19] Angel J. Gallego and Roman Orus. Language design as information renormalization. *arXiv preprint arXiv:1708.01525*, 2019.
- [GPC20] Ivan Glasser, Nicola Pancotti, and J. Ignacio Cirac. From probabilistic graphical models to generalized tensor networks for supervised learning. *IEEE Access*, 8:68169–68182, 2020. <https://doi.org/10.1109/ACCESS.2020.2986279>.
- [Gro15] Misha Gromov. Memorandum Ergo, 2015. Available online: <https://www.ihes.fr/~gromov/wp-content/uploads/2018/08/ergo-cut-copyOct29.pdf>. Accessed on June 1, 2021.
- [Kel82] G.M. Kelly. *Basic Concepts of Enriched Category Theory*. London Mathematical Society Lecture Note Series 64. Cambridge University Press, 1982.
- [KMH⁺20] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [Law73] F. William Lawvere. Metric spaces, generalized logic, and closed categories. *Rendiconti del Seminario Matematico e Fisico di Milano*, 43(1):135–166, 1973. Reprinted in *Reprints in Theory and Applications of Categories* (2002), 1–37. <https://doi.org/10.1007/BF02924844>.
- [Law86] F. William Lawvere. Taking categories seriously. *Revista Colombiana de Matematicas*, XX:147–178, 1986. reprinted as *Reprints in theory and applications of categories*, No. 8 (2005), 1–24.
- [Lei14] Tom Leinster. *Basic Category Theory*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 2014. <https://doi.org/10.1017/CB09781107360068>.
- [LT17] Henry Lin and Max Tegmark. Critical behavior in physics and probabilistic formal languages. *Entropy*, 19(7):299, 2017. <https://doi.org/10.3390/e19070299>.
- [LWM19] Qiuchi Li, Benyou Wang, and Massimo Melucci. CNM: An interpretable complex-valued network for matching. 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 4139–4148, Minneapolis, Minnesota, jun 2019. Association for Computational Linguistics. <https://doi.org/10.18653/v1/N19-1420>.
- [LZSH16] Jingfei Li, Peng Zhang, Dawei Song, and Yuexian Hou. An adaptive contextual quantum language model. *Physica A: Statistical Mechanics and its Applications*, 456:51–67, 2016. <https://doi.org/10.1016/j.physa.2016.03.003>.
- [ML20] Francois Meyer and Martha Lewis. Modelling lexical ambiguity with density matrices. In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 276–290, Online, 2020. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.conll-1.21>.
- [MRT21] Jacob Miller, Guillaume Rabusseau, and John Terilla. Tensor networks for probabilistic sequence modeling. In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 3079–3087. PMLR, 13–15 Apr 2021.
- [MVRL20] John Martyn, Guifre Vidal, Chase Roberts, and Stefan Leichenauer. Entanglement and tensor networks for supervised image classification. *arXiv preprint arXiv:2007.06082*, 2020.
- [Orú14] Román Orús. A practical introduction to tensor networks: Matrix product states

- and projected entangled pair states. *Annals of Physics*, 349:117–158, 2014. <https://doi.org/10.1016/j.aop.2014.06.013>.
- [Orú19] Román Orús. Tensor networks for complex quantum systems. *Nature Reviews Physics*, 1, 08 2019. <https://doi.org/10.1038/s42254-019-0086-7>.
- [Ose11] Ivan Oseledets. Tensor-train decomposition. *SIAM Journal of Scientific Computing*, 33(5):2295–2317, 2011. <https://doi.org/10.1137/090752286>.
- [Pen71] Roger Penrose. Applications of negative dimensional tensors. *Combinatorial mathematics and its applications*, 1:221–244, 1971.
- [PKCS15] Robin Piedeleu, Dimitri Kartsaklis, Bob Coecke, and Mehrnoosh Sadrzadeh. Open system categorical quantum semantics in natural language processing. In Lawrence S. Moss and Pawel Sobocinski, editors, *CALCO*, volume 35 of *LIPICs*, pages 270–289. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2015. <https://doi.org/10.4230/LIPICs.CALCO.2015.270>.
- [PTV17] Vasily Pestun, John Terilla, and Yiannis Vlassopoulos. Language as a matrix product state. *arXiv e-print arXiv:1711.01416*, 2017.
- [PV17] Vasily Pestun and Yiannis Vlassopoulos. Tensor network language model. *arXiv e-print arXiv:1710.10248*, 2017.
- [Rie17] Emily Riehl. *Category Theory in Context*. Dover Publications, 2017.
- [RL19] Chase Roberts and Stefan Leichenauer. Introducing TensorNetwork, an open source library for efficient tensor calculations. *Google AI Blog*, 2019. <https://ai.googleblog.com/2019/06/introducing-tensornetwork-open-source.html>. Accessed on March 1, 2020.
- [RNSS18] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018. Technical report, OpenAI.
- [RS20] Justin Reyes and E. Miles Stoudenmire. A multi-scale tensor network architecture for classification and regression. *Mach. Learn.: Sci. Technol.*, 2:035036, 2020. <https://doi.org/10.1088/2632-2153/abffe8>.
- [Sch11] Ulrich Schollwöck. The density-matrix renormalization group in the age of matrix product states. *Annals of Physics*, 326(1):96–192, 2011. <https://doi.org/10.1016/j.aop.2010.09.012>.
- [SKB18] Mehrnoosh Sadrzadeh, Dimitri Kartsaklis, and Esmā Balkir. Sentence entailment in compositional distributional semantics. *Annals of Mathematics and Artificial Intelligence*, 82(4):189–218, 2018. <https://doi.org/10.1007/s10472-017-9570-x>.
- [SNB13] Alessandro Sordani, Jian-Yun Nie, and Yoshua Bengio. Modeling term dependencies with quantum language models for IR. SIGIR '13, pages 653–662, New York, NY, USA, 2013. Association for Computing Machinery. <https://doi.org/10.1145/2484028.2484098>.
- [SS16] E. Miles Stoudenmire and David J. Schwab. Supervised learning with quantum-inspired tensor networks. *Advances in Neural Information Processing Systems (NIPS)*, 29:4799–4807, 2016.
- [ST19] James Stokes and John Terilla. Probabilistic modeling with matrix product states. *Entropy*, 21(12), 2019. <https://doi.org/10.3390/e21121236>.
- [Sto19] E. Miles Stoudenmire. The tensor network, 2019. Available online: <http://tensornetwork.org>.
- [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, pages 6000–6010, 2017.
- [Wil13] Simon Willerton. Tight spans, Isbell completions and semi-tropical modules. *Theory and Applications of Categories*, 28(22):696–732, 2013.
- [WRVL20] Jinhui Wang, Chase Roberts, Guifre Vidal, and Stefan Leichenauer. Anomaly detection with tensor networks. *arXiv preprint arXiv:2006.02516*, 2020.
- [ZNS⁺18] Peng Zhang, Jiabin Niu, Zhan Su, Benyou Wang, Liqun Ma, and Dawei Song. End-to-end quantum-like language models with application to question answering. In *Proc. 32nd AAAI Conf. Artif. Intell.*, pages 5666–5673, Feb. 2018. AAAI Conference on Artificial Intelligence.

- [ZSZ⁺18] Peng Zhang, Zhan Su, Lipeng Zhang, Benyou Wang, and Dawei Song. A quantum many-body wave function inspired language modeling approach. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM '18*, pages 1303–1312. Association for Computing Machinery, 2018. <https://doi.org/10.1145/3269206.3271723>.
- [ZZM⁺19] Lipeng Zhang, Peng Zhang, Xindian Ma, Shuqin Gu, Zhan Su, and Dawei Song. A generalized language model in tensor space. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):7450–7458, Jul. 2019. <https://doi.org/10.1609/aaai.v33i01.33017450>.