

TOPICS IN COMPUTATIONAL LINGUISTICS  
Neural Network Models of Linguistic Structure  
  
**COURSE NOTES**

Sophie Hao  
Yale University  
LING 380/780

FALL 2021



# Syllabus

## Lecture Schedule (Tentative)

- **September 2.** Introduction, syllabus, course logistics (class outside due to power outage)

### Chapter 1: Semantics

*What is the meaning of meaning? We start the course by learning how modern machine learning algorithms perceive the world through the language of linear algebra.*

- **September 7.** Semantics: features, linear algebra, embeddings
- **September 9.** Semantics: The word2vec model, logistic regression

### Chapter 2: Learning

*We introduce the basic principles of machine learning and the framework of neural networks.*

- **September 14.** Machine learning basics: datasets, model architectures, loss functions, optimization algorithms
- **September 16.** Introduction to neural networks: simple and multi-layer perceptrons, activation functions, linear separability, word2vec as neural networks
- **September 21.** Introduction to neural networks (continued):  $n$ -gram causal language modeling
- **September 23.** Optimization: gradient descent, stochastic gradient descent, gradient computation
- **September 28.** Optimization (continued): Training/development/testing, overfitting, generalization, regularization
- **September 30.** The backpropagation algorithm: computation graphs, automatic differentiation, the linear layer
- **October 5.** The backpropagation algorithm (continued)

- **October 7.** The PyTorch library

## Chapter 3: Recurrent Models

*We learn about specialized neural networks designed to model natural language.*

- **October 12.** Recurrent neural networks: simple recurrent networks, long short-term memory networks, gated recurrent networks
- **October 14.** Recurrent neural networks (continued): causal language modeling
- **October 19.** Syntactic parsing: tree-structured recurrent neural networks, recurrent neural network grammars, shift–reduce parsers
- **October 26.** Sequence-to-sequence tasks: encoder–decoder networks
- **October 28.** Sequence-to-sequence tasks (continued): Bahdanau attention and Luong attention

## Chapter 4: Transformer Models

*We learn about the state of the art in neural network models of language: transfer learning via pre-trained Transformer networks.*

- **November 2.** Transformer networks: self-attention and the Transformer architecture
- **November 4.** Transfer learning: pre-training and fine-tuning, the BERT model

## Journal Club

*Students will become teachers and experience one of the most important activities of computational linguistics researchers: participating in a journal club. In small groups, students will read a research paper from the literature, present what they have learned to the class, and lead a group discussion.*

- **November 9 and 11.** Neural network models of linguistic structure: Do neural networks understand natural language grammar?
- **November 16 and 18.** Neural network models of cognition and language processing: How can AI help us understand ourselves?
- **November 30 and December 2.** Neural network models of language acquisition: inductive bias, poverty of the stimulus arguments
- **December 7 and December 9.** Ethics of language technology: fairness, social bias, environmental impact of deep learning

# Contents

<b>1</b>	<b>Semantics</b>	<b>7</b>
1.1	What Is Meaning?	8
1.2	The Geometry of Words	11
1.2.1	Feature Spaces	14
1.2.2	Embeddings and Dimensionality Reduction	16
1.3	The Distributional Hypothesis	20
1.4	The word2vec Models of Distributional Semantics	21
1.4.1	Binary Classification	22
1.4.2	Logistic Regression	23
1.5	Optional Reading: Dimensionality Reduction	26
1.5.1	Principal Components Analysis	26
1.5.2	t-Distributed Stochastic Neighbors Embedding	32



## CHAPTER 1 Semantics

In linguistics, *semantics* is the study of meaning. In this chapter we will learn how word meanings are represented in neural NLP systems.

### IN-CLASS EXERCISE

#### (1.1) EXERCISE

*Is a burrito a sandwich? Consider that the following items are often referred to as “sandwiches”:*

- a. a peanut butter and jelly sandwich (with two slices of bread)*
- b. a submarine sandwich, also known in Connecticut as a “grinder” (consisting of a single long loaf of bread with a longitudinal slice that contains the fillings)*
- c. a Scandinavian open-faced sandwich (consisting of a single slice of bread with toppings)*
- d. a Gyro sandwich (consisting of meat and vegetables wrapped with a piece of flatbread).*

This question was litigated in the lawsuit *White City Shopping Center v. Panera Bread*.<sup>1</sup> Panera Bread, a fast-casual sandwich restaurant, was given the exclusive right to sell sandwiches at the White City Shopping Center in Shrewsbury, Massachusetts. When Panera found out that White City was negotiating a contract with Qdoba, a fast-casual Mexican restaurant specializing in burritos, they sued the shopping center for breach of contract. Their argument was that burritos are sandwiches, and therefore allowing Qdoba to open a restaurant at the shopping center violated Panera’s exclusive right to sell sandwiches. In other words, whether or not White City was allowed to negotiate a contract with Qdoba crucially depended on the meaning of the word *sandwich*.

#### (1.2) EXERCISE

*What do you think the verdict should be? Did White City Shopping Center breach its contract with Panera by allowing Qdoba to sell burritos?*

The judge’s final opinion states (emphasis added):

Given that the term “sandwiches” is not ambiguous and [Panera’s] Lease does not provide a definition of it, this court applies the **ordinary meaning** of the word. The New Webster Third International Dictionary describes a “sandwich”

as “two thin pieces of bread, usually buttered, with a thin layer (as of meat, cheese, or savory mixture) spread between them.” ... Under this definition and **as dictated by common sense**, this court finds that the term “sandwich” **is not commonly understood to include** burritos, tacos, and quesadillas, which are typically made with a single tortilla and stuffed with a choice filling of meat, rice, and beans. As such, there is no viable legal basis for barring White City from leasing to [Qdoba].

(1.3) EXERCISE

*What is the judge’s rationale? What does it mean for a word to “mean” something?*

## 1.1 What Is Meaning?

What exactly is meaning? Philosophers of language spend a lot of time thinking about this question, and about what certain natural language expressions mean. Linguist Leonard Bloomfield<sup>2</sup> believed that the nature of meaning was too complex for linguists to describe and understand. In an influential book, he wrote:

We can define the meaning of a speech-form accurately when this meaning has to do with some matter of which we possess scientific knowledge ... but we have no precise way of defining words like *love* or *hate*, which concern situations that have not been accurately classified – and these latter are in the great majority.

Moreover, even where we have some scientific (that is, universally recognized and accurate) classification, we often find that the meanings of a language do not agree with this classification. The whale is in German called a ‘fish’: *Walfisch* ... and the bat a ‘mouse’: *Fledermaus* .... Physicists view the color-spectrum as a continuous scale of light-waves of different lengths, ranging from 40 to 72 hundred-thousandths of a millimeter, but languages mark off different parts of this scale quite arbitrarily and without precise limits, in the meanings of such color-names as *violet*, *blue*, *green*, *yellow*, *orange*, *red*, and the color-names of different languages do not embrace the same gradations. The kinship of persons seems a simple matter, but the terminologies of kinship that are used in various languages are extremely hard to analyze.

The statement of meanings is therefore the weak point in language-study, and will remain so until human knowledge advances very far beyond its present state. (Bloomfield, 1933, ch. 9)

According to Bloomfield, each word has a “true meaning,” and it is the job of experts to study those meanings and describe them carefully and precisely. Pluto, for example, would not count as a planet, because astronomers have spent years thinking about what it means to be a “planet,” and decided that Pluto doesn’t fit the definition. But language is full of common-sense concepts like “love” or

<sup>1</sup><https://casetext.com/case/white-city-v-pr-restaurants>

<sup>2</sup>He was a Sterling Professor at Yale from 1940 to 1949!



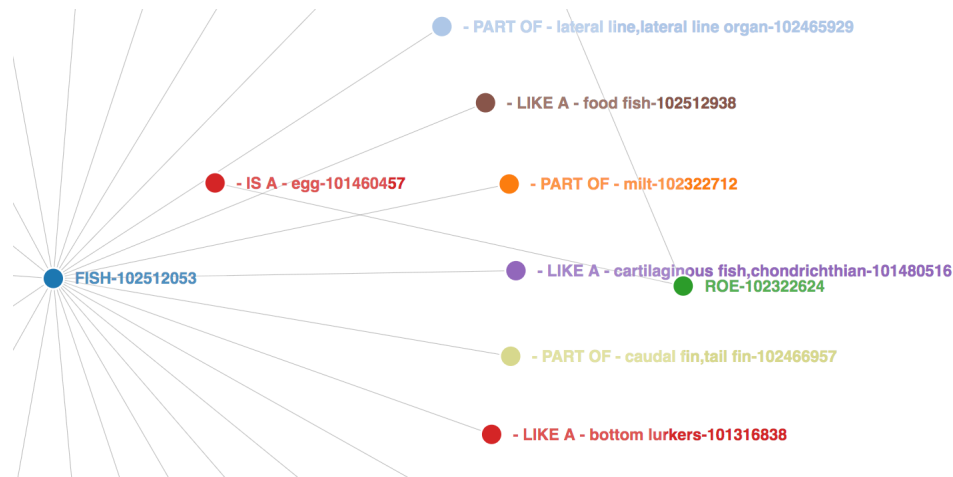


FIGURE 1.1: Relationships between word senses in the WordNet ontology.

“hate” that are difficult to define precisely even if we understand them intuitively. Bloomfield argues in the passage above that we cannot truly understand the nature of meaning until humankind has a better understanding of the concepts that we talk about on a day-to-day basis.

## (1.4) EXERCISE

*How would Bloomfield define a sandwich? Would Bloomfield consider a burrito to be a sandwich? Would Bloomfield agree with the judge’s opinion in *White City Shopping Center v. Panera Bread*?*

## (1.5) EXERCISE

*Do you agree with Bloomfield’s understanding of meaning, and do you think it’s appropriate for language technology? Why or why not?*

In NLP, Bloomfield’s idea of meaning is implemented through *ontologies*. An ontology tries to do what Bloomfield thought was impossible: it is a huge database that contains everything a typical person might know in a machine-readable format. One example of an ontology is WordNet,<sup>3</sup> a database that contains relationships between different word meanings. For each word, the database contains a number of *word senses*, or different meanings of that word. For example, here are the word senses for the word *fish* in WordNet.

(1.6) Word senses for *fish*

- a. (noun) Any of various mostly cold-blooded aquatic vertebrates usually having scales and breathing through gills
- b. (noun) The flesh of fish used as food
- c. (noun, astrology) A person who is born while the sun is in Pisces
- d. (noun) The twelfth sign of the zodiac; the sun is in this sign from about February 19 to March 20

<sup>3</sup><https://wordnet.princeton.edu/>

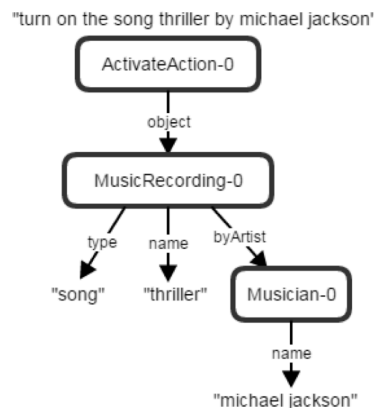


FIGURE 1.2: The abstract representation of the English sentence *turn on the song thriller by michael jackson* used by Amazon Alexa (from Kollar et al., 2018). The rectangles represent ontology objects.

- e. (verb) Seek indirectly
- f. (verb) Catch or try to catch fish or shellfish

The word senses are connected to one another by a graph that represents various relationships between senses. Figure 1.1 shows that according to WordNet, a caudal fin or tail fin is PART OF a fish; a cartilaginous fish or chondrichthian is LIKE A fish; and a roe IS A egg.<sup>4</sup>

While WordNet aims to describe the entire English language, most ontologies employed in NLP systems are limited to a narrow domain of application. For example, Amazon Alexa contains an ontology that describes the various actions (or “skills”) that Alexa can do. A technical paper describes the ontology as follows:<sup>5</sup>

The Alexa ontology ... consists of five primary components:

- **Classes** A hierarchy of Classes, also referred to as types, is defined in the ontology. This hierarchy is a rooted tree, with finer-grained types at deeper levels. Coarse types that are children of THING include PERSON, PLACE, INTANGIBLE, ACTION, PRODUCT, CREATIVEWORK, EVENT and ORGANIZATION. Finer-grained types include MUSICIAN and RESTAURANT.
- **Properties** A given class contains a list of properties, which relate that class to other classes. Properties are defined in a hierarchy, with finer-grained classes inheriting the properties of its parent. There are range restrictions on the available types for both the domain and range of the property.
- **Actions** A hierarchy of actions are defined as classes within the ontology. ACTIONS cover the core functionality of Alexa.
- **Roles** ACTIONS operate on entities via roles. The most common role for an ACTION is the **.object** role, which is defined to be the entity on which the ACTION operates.

<sup>4</sup>You can look at other WordNet graphs at: <https://www.visual-thesaurus.com/wordnet.php>

<sup>5</sup>The documentation is available here: <https://amzn.to/3jrjR9Z>

- **Operators and Relations** A hierarchy of operators and relations represent complex relationships that cannot be expressed easily as properties. Represented as classes, these include ComparativeOperator, Equals and Coordinator ....  
(Kollar et al., 2018)

The ontology objects are used to convert natural-language *requests* into an abstract representation of what the user wants Alexa to do. In Figure 1.2, the request *turn on the song thriller by michael jackson* is interpreted as an **ACTIVATEACTION** whose **.object** is a **MUSICRECORDING**. This **MUSICRECORDING**, in turn, has the **.byArtist** relation with a **MUSICIAN** whose **.name** is “michael jackson.” The various concepts that might appear in a request, like **ACTIVATEACTION**, **MUSICRECORDING**, and **MUSICIAN**, are all defined in the ontology.

## 1.2 The Geometry of Words

Ontologies are an important component of language-based technologies, but they are orthogonal to the material of this course. In this chapter, our goal will be to develop a representation of meaning that is compatible with neural networks. To do so, we will need to be proficient in *linear algebra*, the language of neural networks. Linear algebra is the branch of mathematics that deals with *vectors* and *matrices* and their geometric interpretations. We briefly go over some of the basic definitions here; but because knowledge of linear algebra is a pre-requisite for this course, this should all be review for you, except for the definition of the Hadamard product.

### (1.7) DEFINITION

An  $m \times n$  matrix is a collection of real numbers arranged into a grid with  $m$  rows and  $n$  columns. We denote by  $\mathbb{R}^{m \times n}$  the set of all  $m \times n$  matrices. We use bold, italic, uppercase letters to refer to matrices; e.g.:  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$ . For a matrix  $\mathbf{A}$ , the entry in row  $i$  and column  $j$  is denoted by  $A_{i,j}$ , where  $1 \leq i \leq m$  and  $1 \leq j \leq n$ . We use the following notation to refer to parts of a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ :

- $\mathbf{A}_{i,:}$  is the  $i$ th row of  $\mathbf{A}$  and  $\mathbf{A}_{:,j}$  is the  $j$ th column of  $\mathbf{A}$
- $\mathbf{A}_{i,j:k}$  consists of row  $i$  and columns  $j$  through  $k$  of  $\mathbf{A}$  (inclusive)
- $\mathbf{A}_{i,j:}$  consists of row  $i$  and columns  $j$  through  $n$  of  $\mathbf{A}$  (inclusive)
- $\mathbf{A}_{i,:k}$  consists of row  $i$  and columns 1 through  $k$  of  $\mathbf{A}$  (inclusive)
- $\mathbf{A}_{i:,j}$ ,  $\mathbf{A}_{:k,j}$ , and  $\mathbf{A}_{i:h,j:k}$  are defined analogously for  $1 \leq i < h \leq m$  and  $1 \leq j < k \leq n$ .

A *vector* or *column vector* is a matrix with only one column. We refer to vectors using bold, italic, lowercase letters such as  $\mathbf{a}, \mathbf{b}, \mathbf{c}, \dots$ , and we use the notation  $\mathbb{R}^m$  to refer to the set  $\mathbb{R}^{m \times 1}$ . A *row vector* is a matrix with only one row; i.e., an element of the set  $\mathbb{R}^{1 \times n}$  for some  $n$ . We refer to row vectors using the notation  $\mathbf{a}^\top$ , where  $\mathbf{a}$  is the column vector with the same entries as  $\mathbf{a}^\top$ . The *dimension* of a column vector  $\mathbf{a} \in \mathbb{R}^m$  is  $m$ ; the *dimension* of a row vector  $\mathbf{b}^\top \in \mathbb{R}^{1 \times n}$  is  $n$ . The  $i$ th entry of a column vector  $\mathbf{v}$  is denoted by  $v_i$ . The notations  $\mathbf{v}_{i:}$ ,  $\mathbf{v}_{:,j}$ , and  $\mathbf{v}_{i:j}$  are defined analogously to matrices.

### (1.8) DEFINITION

For matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$  and scalar  $c \in \mathbb{R}$ , the following operations are defined.

- *Matrix addition:*

$$\mathbf{A} + \mathbf{B} = \begin{bmatrix} A_{1,1} + B_{1,1} & A_{1,2} + B_{1,2} & \dots & A_{1,n} + B_{1,n} \\ A_{2,1} + B_{2,1} & A_{2,2} + B_{2,2} & \dots & A_{2,n} + B_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m,1} + B_{m,1} & A_{m,2} + B_{m,2} & \dots & A_{m,n} + B_{m,n} \end{bmatrix}$$

- *Matrix subtraction:*

$$\mathbf{A} - \mathbf{B} = \begin{bmatrix} A_{1,1} - B_{1,1} & A_{1,2} - B_{1,2} & \dots & A_{1,n} - B_{1,n} \\ A_{2,1} - B_{2,1} & A_{2,2} - B_{2,2} & \dots & A_{2,n} - B_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m,1} - B_{m,1} & A_{m,2} - B_{m,2} & \dots & A_{m,n} - B_{m,n} \end{bmatrix}$$

- *Hadamard product:*

$$\mathbf{A} \odot \mathbf{B} = \begin{bmatrix} A_{1,1}B_{1,1} & A_{1,2}B_{1,2} & \dots & A_{1,n}B_{1,n} \\ A_{2,1}B_{2,1} & A_{2,2}B_{2,2} & \dots & A_{2,n}B_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m,1}B_{m,1} & A_{m,2}B_{m,2} & \dots & A_{m,n}B_{m,n} \end{bmatrix}$$

- *Hadamard quotient:*

$$\frac{\mathbf{A}}{\mathbf{B}} = \mathbf{A} / \mathbf{B} = \begin{bmatrix} A_{1,1}/B_{1,1} & A_{1,2}/B_{1,2} & \dots & A_{1,n}/B_{1,n} \\ A_{2,1}/B_{2,1} & A_{2,2}/B_{2,2} & \dots & A_{2,n}/B_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m,1}/B_{m,1} & A_{m,2}/B_{m,2} & \dots & A_{m,n}/B_{m,n} \end{bmatrix}$$

- *Scalar multiplication:*

$$c\mathbf{A} = \mathbf{A}c = \begin{bmatrix} cA_{1,1} & cA_{1,2} & \dots & cA_{1,n} \\ cA_{2,1} & cA_{2,2} & \dots & cA_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ cA_{m,1} & cA_{m,2} & \dots & cA_{m,n} \end{bmatrix}$$

- *Transpose:*

$$\mathbf{A}^\top = \begin{bmatrix} A_{1,1} & A_{2,1} & \dots & A_{m,1} \\ A_{1,2} & A_{2,2} & \dots & A_{m,2} \\ \vdots & \vdots & \ddots & \vdots \\ A_{1,n} & A_{2,n} & \dots & A_{m,n} \end{bmatrix}$$

If  $f : \mathbb{R} \rightarrow \mathbb{R}$  is a function mapping scalars to scalars, we extend  $f$  to matrices by applying them elementwise.

$$f(\mathbf{A}) = \begin{bmatrix} f(A_{1,1}) & f(A_{1,2}) & \dots & f(A_{1,n}) \\ f(A_{2,1}) & f(A_{2,2}) & \dots & f(A_{2,n}) \\ \vdots & \vdots & \ddots & \vdots \\ f(A_{m,1}) & f(A_{m,2}) & \dots & f(A_{m,n}) \end{bmatrix}$$

For vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ , we refer to the operation  $\mathbf{u} + \mathbf{v}$  as *vector addition* instead of matrix addition and  $\mathbf{u} - \mathbf{v}$  as *vector subtraction* instead of matrix subtraction. We define the following operations for vectors.

- *Dot product* (note that this is defined between the row vector  $\mathbf{u}^\top$  and the column vector  $\mathbf{v}$ ):

$$\mathbf{u}^\top \mathbf{v} = \sum_{i=1}^n u_i v_i$$

- *Vector projection* (the projection of  $\mathbf{u}$  onto  $\mathbf{v}$ ):

$$\text{proj}_{\mathbf{v}}(\mathbf{u}) = \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{v}\|^2} \mathbf{v}$$

- *Vector norm* (a.k.a. *vector length*):

$$\|\mathbf{u}\| = \sqrt{\mathbf{u}^\top \mathbf{u}}$$

Finally, when  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{B} \in \mathbb{R}^{n \times p}$ , we define the operation of *matrix multiplication* by

$$\mathbf{AB} = \begin{bmatrix} \mathbf{A}_{1,:} \mathbf{B}_{:,1} & \mathbf{A}_{1,:} \mathbf{B}_{:,2} & \dots & \mathbf{A}_{1,:} \mathbf{B}_{:,p} \\ \mathbf{A}_{2,:} \mathbf{B}_{:,1} & \mathbf{A}_{2,:} \mathbf{B}_{:,2} & \dots & \mathbf{A}_{2,:} \mathbf{B}_{:,p} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{A}_{m,:} \mathbf{B}_{:,1} & \mathbf{A}_{m,:} \mathbf{B}_{:,2} & \dots & \mathbf{A}_{m,:} \mathbf{B}_{:,p} \end{bmatrix}.$$

(1.9) **DEFINITION**

Let  $\mathbb{B} \subseteq \mathbb{R}^n$  be a set of  $n$ -dimensional vectors. The *span* of  $\mathbb{B}$  is the set

$$\text{span}(\mathbb{B}) = \left\{ a_1 \mathbf{b}^{(1)} + a_2 \mathbf{b}^{(2)} + \dots + a_k \mathbf{b}^{(k)} \mid \mathbf{b}^{(1)}, \mathbf{b}^{(2)}, \dots, \mathbf{b}^{(k)} \in \mathbb{B} \text{ and } a_1, a_2, \dots, a_k \in \mathbb{R} \right\}.$$

A set  $\mathbb{A} \subseteq \mathbb{R}^n$  is a *vector space* if  $\mathbb{A} = \text{span}(\mathbb{B})$  for some  $\mathbb{B} \subseteq \mathbb{R}^n$ . We say that  $\mathbb{B}$  is a *basis* for  $\mathbb{A}$  if  $\text{span}(\mathbb{B}) = \mathbb{A}$  and for every set  $\mathbb{C} \subseteq \mathbb{R}^n$ , if  $\text{span}(\mathbb{C}) = \mathbb{A}$ , then  $|\mathbb{C}| \geq |\mathbb{B}|$ . The *dimension* of a vector space  $\mathbb{A}$  is  $\dim(\mathbb{A}) = |\mathbb{B}|$ , where  $\mathbb{B}$  is a basis for  $\mathbb{A}$ .

For reasons that will become apparent in the next chapter, we will use an unusual definition of linear and bilinear maps that allows a constant “bias term” to be added.

(1.10) **DEFINITION**

A *linear map* is a function  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$  of the form

$$f(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$

for some matrix  $\mathbf{W} \in \mathbb{R}^{n \times m}$  and vector  $\mathbf{b} \in \mathbb{R}^n$ . We refer to  $\mathbf{W}$  as a *weight matrix* and  $\mathbf{b}$  as a *bias vector*.

A *bilinear map* is a function  $g : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$  of the form

$$g(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{W}\mathbf{y} + b$$

for some matrix  $\mathbf{W} \in \mathbb{R}^{m \times n}$  and scalar  $b \in \mathbb{R}$ . Again, we refer to  $\mathbf{W}$  as a *weight matrix*, and we refer to  $b$  as a *bias term*.

	Basil	Crushed Italian Tomatoes	Fresh Basil	Fresh Clams	Fresh Mozzarella	Fresh Native Tomatoes	Garlic	Grated Pecorino Romano	Meatball	Mozzarella	Olive Oil	Oregano	Ricotta Cheese
White Clam Pizza	0	0	0	1	0	0	1	1	0	0	1	1	0
Frank Pepe's Original Tomato Pie	0	1	0	0	0	0	0	1	0	0	1	0	0
Margherita	0	1	1	0	1	0	0	1	0	0	1	0	0
Fresh Tomato Pie	1	0	0	0	0	1	1	1	0	1	1	0	0
Meatball & Ricotta	0	1	0	0	0	0	0	1	1	1	1	0	1

FIGURE 1.3: Selected menu items from Frank Pepe's Pizzeria Napoletana.<sup>6</sup>

## (1.11) EXERCISE

Let  $\mathbb{B} \subseteq \mathbb{R}^n$  be a finite set of  $n$ -dimensional vectors. Let  $\mathbf{B} \in \mathbb{R}^{n \times |\mathbb{B}|}$  be a matrix where each column corresponds to an element of  $\mathbb{B}$ ; i.e.,  $\mathbf{b} \in \mathbb{B}$  if and only if  $\mathbf{b} = \mathbf{B}_{:,j}$  for some  $j$ . Prove that  $\mathbf{v} \in \text{span}(\mathbb{B})$  if and only if  $\mathbf{v} = \mathbf{B}\mathbf{a}$  for some  $\mathbf{a} \in \mathbb{R}^{|\mathbb{B}|}$ .

## 1.2.1 Feature Spaces

In order to apply neural networks to language, we need to be able to translate linguistic objects into the language of linear algebra. Vectors are typically interpreted geometrically as points or arrows in a vector space. Linguistic objects like words or sentences, however, are typically thought of as *discrete* objects: there is no such thing as “half of a word” or “the average of two sentences.” Our challenge is to represent discrete linguistic objects using the *continuous* formalism of linear algebra.

In machine learning, discrete objects are represented linear-algebraically by *feature vectors*. To illustrate, take a look at the matrix in Figure 1.3. Let's call this matrix  $\mathbf{P}$ . This matrix is a feature representation of various menu items from Frank Pepe's Pizzeria Napoletana in Wooster Square.<sup>7</sup> Each column represents a pizza topping, and each row represents one of the pizzas on the menu.  $P_{i,j}$  is 1 if pizza  $i$  has topping  $j$ , and  $P_{i,j}$  is 0 otherwise.

## IN-CLASS EXERCISE

## (1.12) EXERCISE

Let  $\mathbf{p}, \mathbf{q} \in \mathbb{R}^{13}$  represent two pizzas, using the same features as in Figure 1.3. What do  $\mathbf{p} + \mathbf{q}$ ,

<sup>6</sup><https://www.pepespizzeria.com/wp-content/uploads/2021/06/Pepe-Menu.T1.0721.pdf>

<sup>7</sup>If you haven't been, I highly recommend it!

$\mathbf{p} \odot \mathbf{q}$ ,  $\max(\mathbf{p}, \mathbf{q})$ ,  $\mathbf{p}^\top \mathbf{q}$ , and  $\mathbf{p} - \mathbf{q}$  represent? Which ones are valid pizzas?

You might think that  $\mathbf{p} - \mathbf{q}$  doesn't represent anything: if  $p_i = 0$  but  $q_i = 1$ , then  $(\mathbf{p} - \mathbf{q})_i = -1$ , and it doesn't make sense for a pizza to have  $-1$  units of topping  $i$ . But vector subtraction does have a useful interpretation – it helps us describe *analogies*.

(1.13) **EXAMPLE**

Define the following pizza vectors.<sup>8</sup>

$$\begin{aligned}\mathbf{p}^\top &= [0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 0] && \text{(white pie with meatballs)} \\ \mathbf{q}^\top &= [0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0] && \text{(plain white pie)} \\ \mathbf{r}^\top &= [0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0] && \text{(red pie with meatballs)} \\ \mathbf{s}^\top &= [0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0] && \text{(plain red pie)}\end{aligned}$$

Observe that  $\mathbf{p} - \mathbf{q} = \mathbf{r} - \mathbf{s}$ . We interpret this equation as representing the following analogy: *a white pie with meatballs is to a plain white pie what a red pie with meatballs is to a plain red pie*. In NLP, we will often rewrite this equation in the form  $\mathbf{p} - \mathbf{q} + \mathbf{s} = \mathbf{r}$ . The vector  $-\mathbf{q} + \mathbf{s}$  then represents the operation of replacing the properties of a plain white pie with those of a plain red pie.

When working with feature vectors, it is important that we be able to measure how similar or different two feature vectors are. The dot product plays a prominent role to that end. During the in-class exercise, you should have realized that the dot product between two pizzas is the number of toppings those two pizzas have in common. Certainly, the more toppings two pizzas share, the more similar they are. However, pizzas with more toppings will have a bigger dot product, and it does not seem fair to say that two pizzas with 50 toppings in common out of 100 are more similar to one another than two pizzas with 4 toppings in common out of 5. To compensate for this phenomenon, we define a measure of similarity that takes the dot product between two feature vectors, and then normalizes this dot product by the vectors' lengths.

(1.14) **DEFINITION**

The *cosine similarity* between two vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$  is defined as

$$\cos(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\| \cdot \|\mathbf{v}\|}.$$

Geometrically,  $\cos(\mathbf{u}, \mathbf{v})$  is the cosine of the angle between the vectors  $\mathbf{u}$  and  $\mathbf{v}$ . This gives us a geometric interpretation of what it means for two feature vectors to share features. When the angle between  $\mathbf{u}$  and  $\mathbf{v}$  is close to 0, and therefore  $\mathbf{u}$  and  $\mathbf{v}$  point in roughly the same direction,  $\cos(\mathbf{u}, \mathbf{v}) \approx 1$ ; and when the angle is roughly  $\pi$ , and therefore  $\mathbf{u}$  and  $\mathbf{v}$  point in roughly opposite directions,  $\cos(\mathbf{u}, \mathbf{v}) \approx -1$ . Thus, when  $\mathbf{u}$  and  $\mathbf{v}$  share a lot of features, they point in roughly the same direction.

<sup>8</sup>In traditional New Haven pizzerias, a *white pie* is a pizza with tomato sauce and mozzarella cheese, while a *red pie* is a pizza with tomato sauce and no cheese.

While the cosine similarity of two vectors  $\mathbf{u}$  and  $\mathbf{v}$  measures how similar  $\mathbf{u}$  and  $\mathbf{v}$  are, we can also use the *cosine distance*, defined as  $1 - \cos(\mathbf{u}, \mathbf{v})$ , to measure how *different* two vectors are.

(1.15) **EXERCISE**

Compute the cosine similarity between every two pizzas in Figure 1.3.

(1.16) **EXERCISE**

Consider the two most dissimilar pizzas possible: the pizza with all the toppings and the pizza with no toppings. What is the angle between their feature vectors?

Now, disregard the pizza example. What does it mean when two feature vectors point in completely opposite directions? Consider the fact that feature values don't have to always be 0 or 1.

## 1.2.2 Embeddings and Dimensionality Reduction

It is hard to derive any intuition from a list of numbers. If we are to understand feature vectors, it would be very helpful if we had a way of visualizing them. Unfortunately, feature spaces typically have lots and lots of features—our pizza vectors have 13 features, but in this course the smallest feature spaces will have at least 300 features—and we can't plot points in spaces of more than 2 dimensions. Therefore, we will need to *embed* our feature vectors into a 2-dimensional space.

In mathematics, an **embedding** is a function that converts objects in one space into analogous objects in another. For example, the circle  $\{\mathbf{x} \in \mathbb{R}^2 \mid \|\mathbf{x}\| = 1\} \subseteq \mathbb{R}^2$  can be embedded in  $\mathbb{R}^3$  as the circle  $\{\mathbf{x} \in \mathbb{R}^3 \mid \|\mathbf{x}\| = 1 \text{ and } x_3 = 0\}$ . Geometrically, we know that the original circle and the embedded circle have the exact same shape because the embedding is *isometric*: the distance between any two points on the original circle is the same as the distance between the analogous two points on the embedded circle.

(1.17) **DEFINITION**

Let  $\mathbb{A}$  and  $\mathbb{B}$  be two arbitrary sets. An *embedding of  $\mathbb{A}$  into  $\mathbb{B}$*  is a one-to-one mapping of  $\mathbb{A}$  to  $\mathbb{B}$ ; i.e., a function  $f : \mathbb{A} \rightarrow \mathbb{B}$  such that for all  $x, y \in \mathbb{A}$ , if  $f(x) = f(y)$ , then  $x = y$ . We say that an embedding  $f : \mathbb{A} \rightarrow \mathbb{B}$  is *isometric with respect to  $d_1$  and  $d_2$*  if  $d_1 : \mathbb{A} \times \mathbb{A} \rightarrow \mathbb{R}$  and  $d_2 : \mathbb{B} \times \mathbb{B} \rightarrow \mathbb{R}$  are functions such that for all  $x, y \in \mathbb{A}$ ,  $d_1(x, y) = d_2(f(x), f(y))$ .

(1.18) **DEFINITION**

The *Euclidean distance* between two vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$  is defined as  $\|\mathbf{u} - \mathbf{v}\|$ .

(1.19) **EXAMPLE**

The function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}^3$  given by  $f(\mathbf{x}) = [x_1 \ x_2 \ 0]^\top$  is an embedding of  $\mathbb{R}^2$  into  $\mathbb{R}^3$ . It is isometric with respect to Euclidean distance and Euclidean distance.

---

<sup>9</sup>This embedding was computed using an algorithm called *multidimensional scaling* (MDS). You don't need to know about MDS, but if you're curious, you can learn about it here: <https://scikit-learn.org/stable/modules/generated/sklearn.manifold.MDS.html>



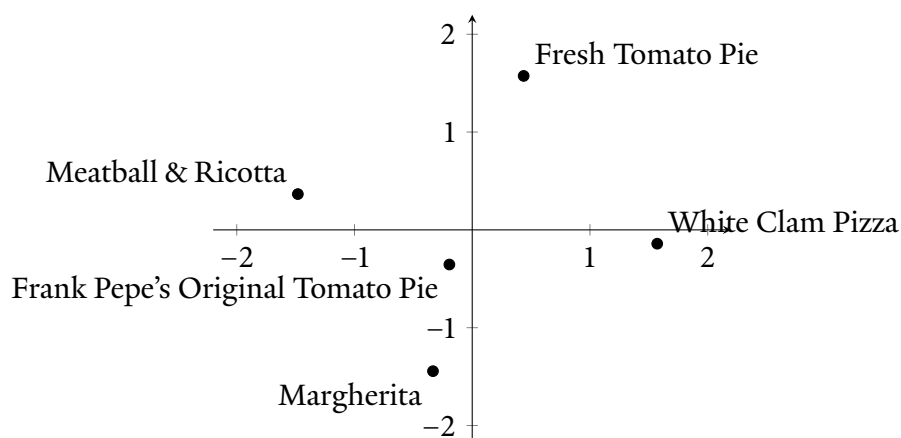


FIGURE 1.4: An embedding of our pizza vectors from  $\mathbb{R}^{13}$  into  $\mathbb{R}^2$ . The embedding is approximately isometric with respect to the absolute value of cosine distance and Euclidean distance.<sup>9</sup>

Figure 1.4 shows an embedding of our pizza vectors into the 2D plane that is approximately isometric with respect to (the absolute value of) cosine distance and Euclidean distance. The closer two pizzas appear to one another within the plot, the greater their cosine similarity. From the diagram, we can quickly observe that Meatball & Ricotta, Frank Pepe’s Original Tomato Pie, and Margherita are relatively similar to one another, while the Fresh Tomato Pie and White Clam Pizza are not so similar. Note that the actual coordinate values don’t mean anything: the axes have no particular interpretation; the only thing that matters is the distance between points.

In data science, the task of embedding points in a high-dimensional space into a lower-dimensional space is known as *dimensionality reduction*. In this course, we will primarily use two algorithms for dimensionality reduction: *principal components analysis* (PCA) and *t-distributed stochastic neighbor embedding* (t-SNE). In the remainder of this section, we’ll go over the basic intuition of how PCA works. A more detailed explanation of both PCA and t-SNE is available in Section 1.5.

The basic idea of PCA is that high-dimensional objects are often well-approximated by lower-dimensional objects, and therefore many of the dimensions can be dropped without losing too much information. For example, look at the data points plotted in Figure 1.5. The points, as a whole, take the form of an elongated ellipse, and a linear model along the major axis of the ellipse would serve as a good approximation of the data. The direction of this line is shown by the longer of the two arrows in the plot, and this direction is chosen in particular because it is the *direction of greatest variance*: the direction that preserves the most information if all points were projected to it, where information is quantified by the variance of the projected data points.

#### (1.20) EXAMPLE

Let’s now use PCA to embed our pizza vectors into  $\mathbb{R}^2$ . We’ll need to use two Python packages: NumPy, a package for doing linear algebra, and scikit-learn, a package that implements basic data science algorithms.

```
1 >>> import numpy as np
2 >>> from sklearn.decomposition import PCA
```

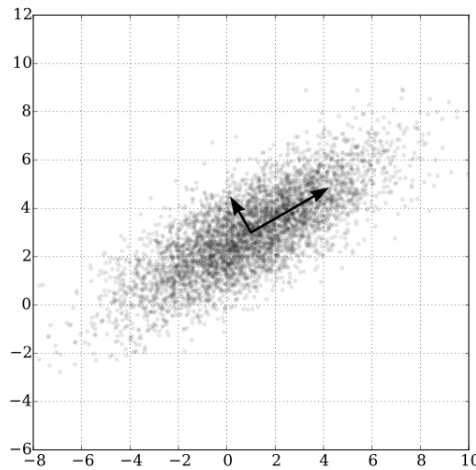


FIGURE 1.5: The two principal components of a 2-dimensional dataset. Image source: Wikimedia Commons user Nicoguaro (CC BY 4.0)

First, we enter in our pizza matrix  $P$  as a NumPy *array*. The array is a data structure that represents vectors and matrices. To create an array, we specify the entries of our matrix by making a list for each row, and then combining the rows into a list of lists.

```
1 >>> p = np.array(
2 ... [[0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0],
3 ... [0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0],
4 ... [0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0],
5 ... [1, 0, 0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0],
6 ... [0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 1]])
```

To perform dimensionality reduction, we create a PCA object and call the `fit_transform` method on our data matrix. We use the `n_components` keyword parameter to specify what vector space we want to embed our vectors into: setting `n_components=2` means we want to embed into  $\mathbb{R}^2$ . The return value of `fit_transform` gives us an array with the embedded coordinates, where each row contains the embedded coordinates of the corresponding row in  $P$ .

```
1 >>> pca = PCA(n_components=2)
2 >>> embedded = pca.fit_transform(p)
3 >>> embedded
4 array([[ 1.01264811, -0.92395855],
5        [-0.50931037, -0.23418494],
6        [-0.9604      , -0.67259834],
7        [ 1.21505519,  0.78969081],
8        [-0.75799292,  1.04105102]])
```

These embeddings are computed by projecting our pizza vectors to the *first two principal components* of our data. The *first principal component* of our data is the direction of greatest

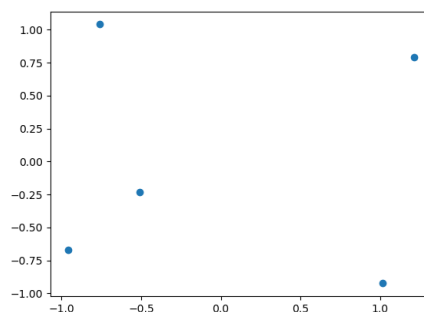


FIGURE 1.6: An embedding of our pizza vectors into  $\mathbb{R}^2$  computed using PCA.

variance, while the *second principal component* is the direction orthogonal to the first principal component with the greatest variance; for  $n > 2$ , the  $n$ th principal component is the direction of greatest variance orthogonal to the first  $n - 1$  principal components. These directions are defined mathematically as unit vectors. You can see the principal components used by the PCA object by looking at the `components_` attribute.

```

1 >>> pca.components_
2 array([[ 0.28534923, -0.52316423, -0.22554482,  0.23781501,
3         -0.22554482,  0.28534923,  0.52316423,  0.
4         -0.17801059,  0.10733863,  0.
5         -0.17801059],
6        [ 0.25736834,  0.04375923, -0.2192067 , -0.30112757,
7         -0.2192067 ,  0.25736834, -0.04375923,  0.
8         0.33928921,  0.59665754,  0.
9         0.33928921]])

```

Finally, let's visualize our embedded vectors. We do this by using the Matplotlib package, which provides tools for creating plots and figures. We first use the `scatter` function to create a scatter plot of our PCA-embedded pizza vectors, and then we use the `savefig` function to save our plot to the file `pca_pizza.png`.

```

1 >>> import matplotlib.pyplot as plt
2 >>> plt.scatter(embedded[:, 0], embedded[:, 1])
3 <matplotlib.collections.PathCollection object at 0x7fdda8245610>
4 >>> plt.savefig("pca_pizza.png")

```

Our PCA plot is shown in Figure 1.6.

#### (1.21) EXERCISE

*What is the Euclidean distance between the PCA embeddings of each pair of pizzas? How do the Euclidean distances compare to the cosine distances between the original 13-dimensional pizza vectors? Is the PCA embedding roughly isometric?*

*Read the documentations for the `distance.euclidean`<sup>10</sup> and `distance.cosine`<sup>11</sup> functions from the `scipy.spatial` package to find your answer.*

## 1.3 The Distributional Hypothesis

Our goal for the rest of this chapter is to figure out how to represent word meanings as feature vectors. Unlike pizzas, however, word meanings cannot be readily reduced to a small collection of features. However hard ontologists might try, the task of parameterizing the entirety of natural language semantics in terms of numerical features is inherently subject to the challenges that Bloomfield discusses in his book. In real life, word meanings are not precisely defined as they are in technical disciplines, and the usage of words varies substantially across individual speakers, communities, demographic groups, and languages.<sup>12</sup> And even if these problems were set aside, developing a comprehensive representation of meaning would require an astronomical amount of engineering labor, presenting a major disadvantage in terms of scalability.

In this course we will assume that it is impossible for engineers to manually define word meanings in terms of feature vectors. Instead, we somehow need to find features for word meanings that can be automatically extracted from data. These features will be *latent*, or *hidden*: they cannot be directly observed in the data, and they do not have any *a priori* interpretation. However, we will make the assumption that word meanings exist in some abstract space where the “similarity” between two meanings can be quantified, and we will desire that our hidden features provide an embedding of word meanings from this abstract space into  $\mathbb{R}^n$  that is isometric with respect to word similarity and cosine similarity. Although the individual hidden features will not have any interpretation, the cosine similarity between feature vectors will be interpreted as the similarity between word meanings, under the assumption that the embedding of words into  $\mathbb{R}^n$  is isometric. Because their interpretations are based on the assumed existence of an isometric embedding, **hidden feature vectors extracted from data that represent word meanings** are known as **word embeddings**.

But how exactly do we find these hidden features? For that, we turn to linguist J. R. Firth (1957), who wrote (with the famous quotation in bold):

As Wittgenstein says, ‘the meaning of words lies in their use.’ The day to day practice of playing language games recognizes customs and rules. It follows that a text in such established usage may contain sentences such as ‘Don’t be such an ass!’, ‘You silly ass!’, ‘What an ass he is!’ In these examples, the word *ass* is in familiar and habitual company, commonly collocated with *you silly* —, *he is a silly* —, *don’t be such an* —. **You shall know a word by the company it keeps!** One of the meanings of *ass* is its habitual collocation with such other words as those above quoted.

<sup>10</sup><https://docs.scipy.org/doc/scipy/reference/generated/scipy.spatial.distance.euclidean.html>

<sup>11</sup><https://docs.scipy.org/doc/scipy/reference/generated/scipy.spatial.distance.cosine.html>

<sup>12</sup>How language might vary is an actively ongoing area of linguistics research, especially here at Yale. Check out the Yale Grammatical Diversity Project for more information: <https://ygdp.yale.edu/>

Bloomfield (1933) “[defines] the *meaning* of a [word] as the situation in which the speaker utters it and the response which it calls forth in the hearer.” According to Bloomfield, we use a certain word, it is because we as speakers want to invoke a certain idea in our listeners; and ontologies assume that the right way to define a word’s meaning is to describe the idea we want to invoke when we use that word. Firth, however, disagrees with Bloomfield. According to Firth, the meaning of a word is entirely captured by how we use that word. Specifically, the “usage” of a word is described by its *collocations*—statistical descriptions of what other words are used together with that word. In the quotation above, the vulgar meaning of *ass* is described by the fact that *ass* is frequently used with phrases like *he is a silly* and *don’t be such an* that are commonly used with with profanity.

Firth’s notion of meaning is known as the *distributional hypothesis*, and it has given rise to methods of creating word embeddings by modeling collocations. Although there are a variety of distributional methods for creating word embeddings, in the lecture-based portion of this course we will focus on just one such method.

#### IN-CLASS EXERCISE

##### (1.22) EXERCISE

*Would Firth consider a burrito to be a sandwich? Why or why not?*

*Do you agree more with Firth’s concept of meaning, or Bloomfield’s? Which concept do you think the verdict from *White City Shopping Center v. Panera Bread* was based on?*

It is interesting to note that the judge’s opinion made reference to both Bloomfield’s and Firth’s notions of meaning: it simultaneously cited a dictionary definition of *sandwich* while also appealing to “common sense.”

## 1.4 The word2vec Models of Distributional Semantics

The *word2vec* models (Mikolov et al., 2013a) are a family of distributional algorithms used to create word embeddings. Not only have the word2vec models been extremely influential in the development of neural NLP, they are also one of the simplest examples of how neural networks can be applied to the modeling of natural language. In the end of this chapter we will learn about one of the word2vec models, the *Skip-Gram with Negative Sampling* model (SGNS, Mikolov et al., 2013b; Goldberg and Levy, 2014), without the framework of neural networks. After we formally introduce neural networks in the next chapter, we will return to word2vec and see how the other models can be understood as neural networks.

The basic idea behind SGNS is to come up with word embeddings such that words that appear together have a high cosine similarity, while words that do not appear together have a low cosine similarity. More formally, we assume that there is a set of words  $\mathbb{V}$  called the *vocabulary*, and to each word  $w \in \mathbb{V}$  we assign a *word embedding*  $\llbracket w \rrbracket \in \mathbb{R}^d$  and a *context embedding*  $\langle w \rangle \in \mathbb{R}^d$ . These embeddings are arranged into matrices  $\mathbf{W}, \mathbf{C} \in \mathbb{R}^{|\mathbb{V}| \times d}$  such that each row of  $\mathbf{W}$  is the word embedding of some word and each row of  $\mathbf{C}$  is the context embedding of some word. The entries of  $\mathbf{W}$  and  $\mathbf{C}$  are *parameters* of the model: we initialize them to random values, and a learning algorithm will adjust

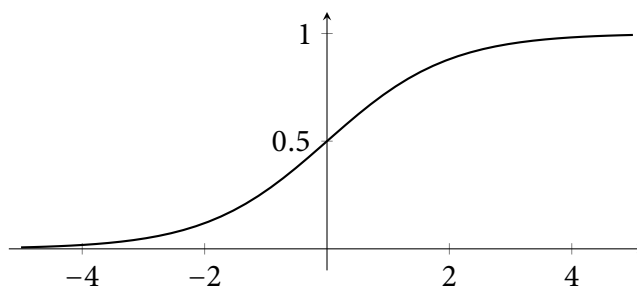


FIGURE 1.7: The sigmoid function, used in logistic regression.

these values to better reflect the data. The constant  $d$ , the dimensionality of our word embedding and context embedding vectors, is a *hyperparameter*: its value needs to be manually defined beforehand, and often times engineers and researchers will try a number of different hyperparameter values in order to see which settings work the best.

### 1.4.1 Binary Classification

In order to operationalize the SGNS model, we conceptualize the task of finding word embeddings as a *binary classification task*. A classification task is any task where feature vectors need to be assigned to one of two or more *classes*; a binary classification task is a classification task with exactly two possible classes. For SGNS, our binary classifier will take a word embedding and a context embedding, and classify this pair of vectors as “occurring together” or “not occurring together.” We will build our classifier by fitting a *logistic model* to a dataset.

(1.23) **DEFINITION**

A *logistic model* is a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  of the form

$$f(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + b)$$

for some  $\mathbf{w} \in \mathbb{R}^n$  and  $b \in \mathbb{R}$ , where  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is the *sigmoid function* given by

$$\sigma(x) = \frac{1}{1 + e^{-x}}.$$

Note that for all  $x \in \mathbb{R}$ ,  $0 < \sigma(x) < 1$ .

The heart of the logistic model is the *sigmoid curve*, shown in Figure 1.7. The sigmoid curve is an S-shaped curve that maps real numbers to numbers between 0 and 1. The sigmoid’s  $y$ -intercept is  $1/2$ , and it is roughly linear near the  $y$  axis. When  $x$  has a large magnitude, however, the sigmoid approaches 1 as  $x$  tends to  $\infty$  and 0 as  $x$  tends to  $-\infty$ . A logistic model combines the sigmoid curve with a linear map. Suppose we are classifying inputs into two classes,  $c_1$  and  $c_2$ . The linear map computes a score that measures the degree to which the input resembles elements of one class or the other. If the score is positive, then the input is more likely to belong to class  $c_1$ , and if the score is negative, then the input is more likely to belong to class  $c_2$ . The sigmoid curve turns this score

$w$	$c$	$y$
dog	cat	1
dog	the	1
dog	big	1
dog	acetic	0
dog	pernicious	0
dog	spherification	0
$\vdots$	$\vdots$	$\vdots$

FIGURE 1.8: An example of what a dataset for the SGNS classification task might look like.

into a number between 0 and 1, which we interpret as the probability that the input belongs to class  $c_1$ .

Now, let's design a logistic model for our word embeddings. Remember that we want  $\cos(\llbracket w \rrbracket, \llbracket c \rrbracket)$  to be high when  $w$  and  $c$  occur together, and low when  $w$  and  $c$  do not occur together. Since cosine similarity is really just a normalized version of the dot product, we can easily turn the cosine similarity into a linear map by dropping the denominator and keeping the dot product. We then combine the dot product with the sigmoid curve to form a logistic model.

## MODEL

(1.24) **Logistic Model for SGNS**

$$\mathbb{P}[\text{words } w \text{ and } c \text{ occur together}] = \sigma(\llbracket w \rrbracket^\top \llbracket c \rrbracket)$$

## 1.4.2 Logistic Regression

Now that we have our logistic model, we need to *fit* the model to a dataset in order to obtain our word and context embedding matrices. The dataset will be a table  $\mathbb{D} \subseteq \mathbb{V} \times \mathbb{V} \times \{0, 1\}$ , with each element of  $\mathbb{D}$  of the form  $(w, c, y)$ , where  $y = 1$  if  $w$  and  $c$  occur together and  $y = 0$  if  $w$  and  $c$  do not occur together. Because the dataset indicates how each data point should be labeled by our logistic classifier, the classification task we are doing is an example of *supervised learning*. The process of fitting a logistic model to a dataset is called *logistic regression*.

To obtain the dataset  $\mathbb{D}$ , we take a *corpus* of text and extract *skip-grams* from it. A skip-gram is a sequence of words, with one word missing from the middle. For example, Figure 1.9 illustrates how the skip-gram of *New York said on \_\_\_\_\_ that he would* is extracted from a New York Times article. The missing word is flanked by 5 words to the left and 3 words to the right; the larger of these two numbers is the skip-gram's *window size*, which in this case is 5. For each word  $w$  in the corpus, we extract a skip-gram around that word, and for each  $c$  in the skip-gram, we add the tuple  $(w, c, 1)$  to  $\mathbb{D}$ . The window size of the skip-gram is chosen randomly, with the maximum possible window size determined by a hyperparameter  $m$ . In order to add examples with a label of  $y = 0$ , we also randomly choose words  $c_1, c_2, \dots, c_{k_n}$  from  $\mathbb{V}$ , where  $k$  is another hyperparameter,  $n \leq 2m$

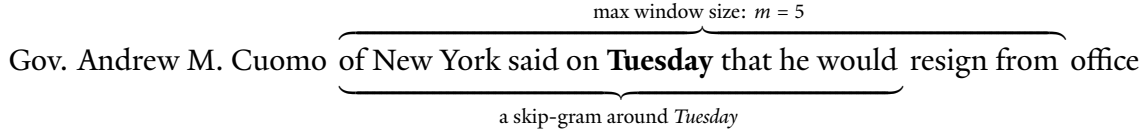


FIGURE 1.9: An example of a skip-gram with window size 5 around the word *Tuesday*.<sup>13</sup>

is the number of words in the skip-gram (excluding the missing word), and for each  $i \leq kn$  we add the items  $(w, c_i, 0)$  to  $\mathbb{D}$ .

After we have obtained our dataset  $\mathbb{D}$ , we solve the following *optimization problem*:

$$\max_{W, C} \left( \prod_{(w, c, 1) \in \mathbb{D}} \mathbb{P}[w \text{ and } c \text{ occur together}] \right) \left( \prod_{(w, c, 0) \in \mathbb{D}} \mathbb{P}[w \text{ and } c \text{ do not occur together}] \right)$$

The problem above is read as follows: “find the values of  $W$  and  $C$  that result in the highest possible value of the expression.” If we assume that each data point of  $\mathbb{D}$  is independent and identically distributed (an assumption that might not be reasonable, but that we will make anyway just because it’s convenient), then the expression being maximized is the probability that a dataset with the same  $(w, c)$  pairs as  $\mathbb{D}$  would be labeled in the same way that  $\mathbb{D}$  is labeled. In other words, our embedding matrices must cause our logistic model to predict that the words that actually do occur together in the real world have a high probability of doing so, while words that do not actually occur together in the real world have a low probability of occurring together. This learning strategy, in which we adjust our model parameters to maximize the probability of our dataset, is called *maximum likelihood estimation*.

Now, we make some simplifications to our optimization problem. First, plugging in our logistic model and observing that  $1 - \sigma(x) = \sigma(-x)$  for all  $x$ , we rewrite the above as follows.

$$\max_{W, C} \left( \prod_{(w, c, 1) \in \mathbb{D}} \sigma(\llbracket w \rrbracket^\top \llbracket c \rrbracket) \right) \left( \prod_{(w, c, 0) \in \mathbb{D}} \sigma(-\llbracket w \rrbracket^\top \llbracket c \rrbracket) \right)$$

Although our expression is now much simpler than it was before, it still has a subtle flaw: it is *numerically unstable*, meaning that it is prone to rounding errors when implemented in code. This is because products of probabilities become very, very small as more and more probabilities are multiplied together. Eventually, the product may accidentally get rounded to 0, at which point the entire expression is zeroed out. To alleviate this problem, we simply take the natural log of the entire expression. Since the natural log is an *increasing function* (i.e., for all  $x, y > 0$  we have  $x > y$  if and only if  $\ln(x) > \ln(y)$ ), maximizing the natural log of an expression is the same thing as maximizing the original expression.

$$\max_{W, C} \left( \sum_{(w, c, 1) \in \mathbb{D}} \ln(\sigma(\llbracket w \rrbracket^\top \llbracket c \rrbracket)) \right) + \left( \sum_{(w, c, 0) \in \mathbb{D}} \ln(\sigma(-\llbracket w \rrbracket^\top \llbracket c \rrbracket)) \right)$$

<sup>13</sup><https://www.nytimes.com/2021/08/10/nyregion/andrew-cuomo-resigns.html>



Next, we multiply our expression by  $-1$  and restate our optimization problem as a *minimization* problem instead of a maximization problem. There is no mathematical, statistical, or computational reason for doing this; we do it simply because regression problems such as linear regression or logistic regression are traditionally thought of as minimization problems.

$$\min_{W, C} - \left( \sum_{(w, c, 1) \in \mathbb{D}} \ln(\sigma(\llbracket w \rrbracket^\top \llbracket c \rrbracket)) \right) - \left( \sum_{(w, c, 0) \in \mathbb{D}} \ln(\sigma(-\llbracket w \rrbracket^\top \llbracket c \rrbracket)) \right)$$

Finally, because our dataset contains  $k$  times as many examples with a label of  $y = 0$  as examples with a label of  $y = 1$ , we scale the terms associated with the label of  $y = 0$  by  $1/k$ . Our final minimization problem is as follows.

$$\min_{W, C} - \left( \sum_{(w, c, 1) \in \mathbb{D}} \ln(\sigma(\llbracket w \rrbracket^\top \llbracket c \rrbracket)) \right) - \frac{1}{k} \left( \sum_{(w, c, 0) \in \mathbb{D}} \ln(\sigma(-\llbracket w \rrbracket^\top \llbracket c \rrbracket)) \right)$$

In the next chapter, we will learn how this minimization problem is solved. Once it is solved, we discard  $C$  and use  $W$  as our word embeddings.

To summarize, the full SGNS model is as follows.

#### MODEL: SKIP-GRAM WITH NEGATIVE SAMPLING

##### Hyperparameters:

- $\mathbb{V}$ , a vocabulary
- $d$ , the dimensionality of the word and context embeddings
- $m$ , the maximum window size of skip-grams
- $k$ , the number of  $y = 0$  examples to include per  $y = 1$  example

##### Parameters:

- $W \in \mathbb{R}^{|\mathbb{V}| \times d}$ , a matrix of word embeddings
- $C \in \mathbb{R}^{|\mathbb{V}| \times d}$ , a matrix of context embeddings

##### Data:

- A corpus of text

##### Procedure:

- Initialize  $\mathbb{D}$  to an empty set.
- For each word  $w$  in the corpus:
  - Let  $c_1, c_2, \dots, c_n$  be a skip-gram around  $w$  with window size at most  $m$ . For each  $i \leq n$ , add  $(w, c_i, 1)$  to  $\mathbb{D}$ .
  - Let  $c'_1, c'_2, \dots, c'_{kn}$  be randomly sampled from  $\mathbb{V}$ . For each  $i \leq kn$ , add  $(w, c'_i, 0)$  to  $\mathbb{D}$ .
- Solve the following minimization problem:

$$W, C \leftarrow \operatorname{argmin}_{W, C} - \left( \sum_{(w, c, 1) \in \mathbb{D}} \ln(\sigma(\llbracket w \rrbracket^\top \llbracket c \rrbracket)) \right) - \frac{1}{k} \left( \sum_{(w, c, 0) \in \mathbb{D}} \ln(\sigma(-\llbracket w \rrbracket^\top \llbracket c \rrbracket)) \right)$$

- Return  $W$ .

## 1.5 Optional Reading: Dimensionality Reduction

In this optional section, we learn about the PCA and  $t$ -SNE algorithms.

### 1.5.1 Principal Components Analysis

We have already seen the basic intuition behind PCA: we want to find the *principal components* of a set of feature vectors, which are defined as the *directions of greatest variance* for those vectors. More concretely, suppose we have a dataset  $X \in \mathbb{R}^{m \times n}$ , represented as a matrix in which each row contains an  $n$ -dimensional feature vector. We want to find a set of vectors

$$\mathbb{B} = \{\psi^{(1)}, \psi^{(2)}, \dots, \psi^{(n)}\}$$

with the following properties.

(1.25) **Axioms for Principal Components**

- a.  $\mathbb{B}$  forms a basis for  $\mathbb{R}^n$ .
- b. Each element of  $\mathbb{B}$  has unit length: for all  $i$ ,  $\|\psi^{(i)}\| = 1$ .
- c. The elements of  $\mathbb{B}$  are pairwise orthogonal: for all  $i$  and  $j$  where  $i \neq j$ ,  $(\psi^{(i)})^\top (\psi^{(j)}) = 0$ . (Geometrically, this means that  $\psi^{(i)}$  is perpendicular to  $\psi^{(j)}$  in  $\mathbb{R}^n$ .)
- d. When  $X$  is re-written using the coordinates of  $\mathbb{B}$ , the  $\psi^{(1)}$  direction contains the greatest possible variation of values; and for  $i > 1$ , the  $\psi^{(i)}$  direction contains the greatest possible variation of values subject to the constraint that  $\psi^{(i)}$  must be orthogonal to  $\psi^{(j)}$  for all  $j < i$ .

The elements of  $\mathbb{B}$  are the principal components of  $X$ ; for each  $i$ ,  $\psi^{(i)}$  is the  $i$ th *principal component* of  $X$ . When conditions (1.25a–c) are satisfied, we refer to  $\mathbb{B}$  as an *orthonormal basis* for  $\mathbb{R}^n$ . Axiom (1.25d), for now, is stated in a hand-wavy manner. To define the PCA algorithm, we will have to first define the concept of “variance,” and then figure out how to find vectors that maximize variance.

**The Covariance Matrix.** The first step of PCA is to define a *covariance matrix* that measures the variance represented by the data of  $X$ . The notion of “variance” we use here is based on statistical variance, from probability theory. Recall that the variance of a scalar-valued random variable is defined as follows.

(1.26) **DEFINITION**

Let  $a$  and  $b$  be scalar-valued random variables. The *covariance* of  $a$  and  $b$  is the quantity

$$\text{cov}(a, b) = \mathbb{E}[(a - \mathbb{E}[a])(b - \mathbb{E}[b])].$$

The *variance* of  $a$  is defined as  $\text{var}(a) = \text{cov}(a, a)$ .

We extend the notion of covariance to vector-valued random variables by treating each position of each vector as a separate variable. The covariance is defined as a matrix containing the covariance between each pair of vector entries.

## (1.27) DEFINITION

Let  $\mathbf{a}$  and  $\mathbf{b}$  be  $n$ -dimensional vector-valued random variables. The *covariance of  $\mathbf{a}$  and  $\mathbf{b}$*  is the matrix

$$\text{cov}(\mathbf{a}, \mathbf{b}) = \mathbb{E}[(\mathbf{a} - \mathbb{E}[\mathbf{a}])(\mathbf{b} - \mathbb{E}[\mathbf{b}])^\top] \in \mathbb{R}^{n \times n}.$$

Notice that  $\text{cov}(\mathbf{a}, \mathbf{b})_{i,j} = \text{cov}(a_i, b_j)$  for all positions  $i$  and  $j$ .

Since the covariance between the two random variables is a matrix, there is no scalar value that represents the “variance” of a vector-valued random variable. Instead, the vector-valued analogue of variance is simply the covariance of a random variable  $\mathbf{a}$  with itself:  $\text{cov}(\mathbf{a}, \mathbf{a})$ .

The covariance matrix for  $\mathbf{X}$  is defined by pretending that  $\mathbf{X}$  is a data sample that estimates the distribution of some  $n$ -dimensional random variable  $\mathbf{x}$ . The expected value of  $\mathbf{x}$  is estimated using the mean of the sample, denoted by  $\bar{\mathbf{X}}$ :

$$\mathbb{E}[\mathbf{x}] = \bar{\mathbf{X}}^\top = \frac{1}{m} \sum_{i=1}^m (\mathbf{X}_{i,:})^\top = \frac{\mathbf{X}^\top \mathbf{1}}{m}.$$

We then estimate  $\text{cov}(\mathbf{x}, \mathbf{x})$  by plugging in the above expression into the formula for covariance:

$$\text{cov}(\mathbf{x}, \mathbf{x}) = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^\top] = \mathbb{E}[(\mathbf{x} - \bar{\mathbf{X}}^\top)(\mathbf{x} - \bar{\mathbf{X}}^\top)^\top]$$

and then using the sample mean of  $(\mathbf{x} - \bar{\mathbf{X}}^\top)(\mathbf{x} - \bar{\mathbf{X}}^\top)^\top$  to estimate the expected value:

$$\mathbb{E}[(\mathbf{x} - \bar{\mathbf{X}}^\top)(\mathbf{x} - \bar{\mathbf{X}}^\top)^\top] = \frac{1}{m} \sum_{i=1}^m ((\mathbf{X}_{i,:})^\top - \bar{\mathbf{X}}^\top)(\mathbf{X}_{i,:} - \bar{\mathbf{X}}) = \frac{(\mathbf{X} - \bar{\mathbf{X}})^\top (\mathbf{X} - \bar{\mathbf{X}})}{m}.$$

We now use this formula to define the covariance matrix.

## (1.28) DEFINITION

Let  $\mathbf{X} \in \mathbb{R}^{m \times n}$  be a matrix. The *covariance matrix of  $\mathbf{X}$*  is the  $n \times n$  matrix

$$\text{cov}(\mathbf{X}) = \frac{(\mathbf{X} - \bar{\mathbf{X}})^\top (\mathbf{X} - \bar{\mathbf{X}})}{m},$$

where  $\bar{\mathbf{X}} = \mathbf{1}^\top \mathbf{X} / m$ .

## (1.29) EXERCISE

*Prove that the last part of the derivation is valid: using the definition of matrix multiplication, convince yourself that*

$$\frac{1}{m} \sum_{i=1}^m ((\mathbf{X}_{i,:})^\top - \bar{\mathbf{X}}^\top)(\mathbf{X}_{i,:} - \bar{\mathbf{X}}) = \frac{(\mathbf{X} - \bar{\mathbf{X}})^\top (\mathbf{X} - \bar{\mathbf{X}})}{m}.$$

**Projecting onto a Basis.** When we find our orthonormal basis  $\mathbb{B}$  containing the principal components of  $\mathbf{X}$ , we will want to re-write the rows of  $\mathbf{X}$  in terms of coordinates with respect to  $\mathbb{B}$ . The coordinates of a vector  $\mathbf{v} \in \mathbb{R}^n$  with respect to  $\mathbb{B}$  are given by a vector  $\mathbf{a} \in \mathbb{R}^n$  such that

$$\mathbf{v} = a_1\psi^{(1)} + a_2\psi^{(2)} + \cdots + a_n\psi^{(n)} = \begin{bmatrix} \left| \psi^{(1)} \right\rangle & \left| \psi^{(2)} \right\rangle & \cdots & \left| \psi^{(n)} \right\rangle \end{bmatrix} \mathbf{a}.$$

By solving for  $\mathbf{a}$ , we obtain the following formula for re-writing  $\mathbf{v}$  in terms of  $\mathbb{B}$ :

$$\mathbf{a} = \begin{bmatrix} \left| \psi^{(1)} \right\rangle & \left| \psi^{(2)} \right\rangle & \cdots & \left| \psi^{(n)} \right\rangle \end{bmatrix}^{-1} \mathbf{v};$$

in other words, to find  $\mathbf{a}$ , we simply multiply  $\mathbf{v}$  on the left by the *change of basis matrix* for  $\mathbb{B}$ .

(1.30) **DEFINITION**

Let

$$\mathbb{B} = \{\mathbf{b}^{(1)}, \mathbf{b}^{(2)}, \dots, \mathbf{b}^{(n)}\}$$

be a basis for  $\mathbb{R}^n$ . The *change of basis matrix* for  $\mathbb{B}$  is the matrix

$$\mathbf{B} = \begin{bmatrix} \left| \mathbf{b}^{(1)} \right\rangle & \left| \mathbf{b}^{(2)} \right\rangle & \cdots & \left| \mathbf{b}^{(n)} \right\rangle \end{bmatrix}^{-1}.$$

For a vector  $\mathbf{v} \in \mathbb{R}^n$ , the vector  $\mathbf{B}\mathbf{v}$  is the *coordinates of  $\mathbf{v}$  with respect to  $\mathbb{B}$* .

Notice that computing the change of basis matrix requires inverting a matrix. This is a nasty operation: inverses don't always exist, and when they do, they can be computationally expensive to find. Thankfully, however, our assumption that  $\mathbb{B}$  is orthonormal will make the inverse operation easy to compute.

(1.31) **THEOREM**

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a square matrix. Then,

- a. each column/row of  $\mathbf{A}$  has unit length and
  - b. the columns/rows of  $\mathbf{A}$  are pairwise orthogonal
- if and only if  $\mathbf{A}^{-1} = \mathbf{A}^\top$ .

*Proof.* Suppose conditions a and b are satisfied by the rows of  $\mathbf{A}$ . Then, observe that

(1.32)

$$(\mathbf{A}\mathbf{A}^\top)_{i,j} = \mathbf{A}_{i,:}(\mathbf{A}_{j,:})^\top = \begin{cases} 1, & i = j \\ 0, & i \neq j, \end{cases}$$

thus  $\mathbf{A}\mathbf{A}^\top = \mathbf{I}_n$ , so  $\mathbf{A}^\top = \mathbf{A}^{-1}$ . Similarly, if a and b are satisfied by the columns of  $\mathbf{A}$ , then

(1.33)

$$(\mathbf{A}^\top \mathbf{A})_{i,j} = (\mathbf{A}_{:,i})^\top \mathbf{A}_{:,j} = \begin{cases} 1, & i = j \\ 0, & i \neq j, \end{cases}$$

so again  $\mathbf{A}\mathbf{A}^\top = \mathbf{I}_n$ , thus  $\mathbf{A}^\top = \mathbf{A}^{-1}$ .

Now, suppose  $\mathbf{A}^{-1} = \mathbf{A}^\top$ . Then,  $\mathbf{I}_n = \mathbf{A}\mathbf{A}^\top = \mathbf{A}^\top \mathbf{A}$ , so equations (1.32) and (1.33) are both satisfied, which in turn implies that conditions a and b are satisfied by both the rows and columns of  $\mathbf{A}$ .  $\square$

Based on this property, the change of basis matrix for  $\mathbb{B}$  is simply

$$\mathbf{B} = \begin{bmatrix} | & | & & | \\ \psi^{(1)} & \psi^{(2)} & \dots & \psi^{(n)} \\ | & | & & | \end{bmatrix}^{-1} = \begin{bmatrix} | & | & & | \\ \psi^{(1)} & \psi^{(2)} & \dots & \psi^{(n)} \\ | & | & & | \end{bmatrix}^\top = \begin{bmatrix} - & \psi^{(1)} & - \\ - & \psi^{(2)} & - \\ - & \vdots & - \\ - & \psi^{(n)} & - \end{bmatrix}.$$

**The Direction of Greatest Variance.** When we re-write  $\mathbf{X}$  using  $\mathbb{B}$ 's coordinates, the resulting data matrix is  $(\mathbf{B}\mathbf{X}^\top)^\top = \mathbf{X}\mathbf{B}^\top$ . For each  $i$ , the  $m$ -dimensional vector  $\mathbf{X}\psi^{(i)}$  contains the component in the  $\psi^{(i)}$  direction for each item of  $\mathbf{X}$ .

Observe that since  $\mathbf{X}\mathbf{v}$  is a column vector as long as  $\mathbf{v} \in \mathbb{R}^n$ , the covariance matrix  $\text{cov}(\mathbf{X}\mathbf{v})$  contains only one row and one column. We can therefore treat  $\text{cov}(\mathbf{X}\mathbf{v})$  as a scalar value, allowing us to define  $\psi^{(1)}$  as the vector  $\mathbf{v}$  that maximizes  $\text{cov}(\mathbf{X}\mathbf{v})$  subject to the constraint that  $\mathbf{v} \neq \mathbf{0}$ :

$$\psi^{(1)} = \underset{\mathbf{v}}{\text{argmax}} \text{cov}(\mathbf{X}\mathbf{v}) \text{ subject to } \mathbf{v} \neq \mathbf{0}.$$

For the other principal components, we do the same thing, but add in a constraint that each principal component must be orthogonal to the previous principal components:

$$\psi^{(i)} = \underset{\mathbf{v}}{\text{argmax}} \text{cov}(\mathbf{X}\mathbf{v}) \text{ subject to } \mathbf{v} \neq \mathbf{0} \text{ and } \forall j < i [\mathbf{v}^\top \psi^{(j)} = 0].$$

Our goal is now to solve these optimization problems.

**The Spectral Theorem.** Let's begin by doing some algebra to figure out the relationship between  $\text{cov}(\mathbf{X}\mathbf{v})$  and  $\text{cov}(\mathbf{X})$ . First, we calculate the mean  $\overline{\mathbf{X}\mathbf{v}}$  as follows:

$$\overline{\mathbf{X}\mathbf{v}} = \frac{\mathbf{1}^\top (\mathbf{X}\mathbf{v})}{m} = \frac{(\mathbf{1}^\top \mathbf{X})\mathbf{v}}{m} = \overline{\mathbf{X}}\mathbf{v}.$$

Plugging this into the definition of the covariance matrix, we get

$$\begin{aligned}
 \text{cov}(\mathbf{X}\mathbf{v}) &= \frac{(\mathbf{X}\mathbf{v} - \overline{\mathbf{X}\mathbf{v}})^\top (\mathbf{X}\mathbf{v} - \overline{\mathbf{X}\mathbf{v}})}{m} \\
 &= \frac{(\mathbf{X}\mathbf{v} - \overline{\mathbf{X}}\mathbf{v})^\top (\mathbf{X}\mathbf{v} - \overline{\mathbf{X}}\mathbf{v})}{m} \\
 &= \frac{(\mathbf{v}^\top \mathbf{X}^\top - \mathbf{v}^\top \overline{\mathbf{X}}^\top)(\mathbf{X}\mathbf{v} - \overline{\mathbf{X}}\mathbf{v})}{m} \\
 &= \frac{\mathbf{v}^\top (\mathbf{X}^\top - \overline{\mathbf{X}}^\top)(\mathbf{X} - \overline{\mathbf{X}})\mathbf{v}}{m} \\
 &= \frac{\mathbf{v}^\top (\mathbf{X} - \overline{\mathbf{X}})^\top (\mathbf{X} - \overline{\mathbf{X}})\mathbf{v}}{m} \\
 &= \mathbf{v}^\top \text{cov}(\mathbf{X})\mathbf{v}.
 \end{aligned}$$

Now, notice that  $\text{cov}(\mathbf{X})$  is *symmetric*; that is to say,  $\text{cov}(\mathbf{X}) = \text{cov}(\mathbf{X})^\top$ :

$$\text{cov}(\mathbf{X})^\top = \left( \frac{(\mathbf{X} - \overline{\mathbf{X}})^\top (\mathbf{X} - \overline{\mathbf{X}})}{m} \right)^\top = \frac{(\mathbf{X} - \overline{\mathbf{X}})^\top ((\mathbf{X} - \overline{\mathbf{X}})^\top)^\top}{m} = \frac{(\mathbf{X} - \overline{\mathbf{X}})^\top (\mathbf{X} - \overline{\mathbf{X}})}{m} = \text{cov}(\mathbf{X}).$$

The *Spectral Theorem* says that when a matrix  $\mathbf{A}$  is symmetric, the maximum value of  $\mathbf{v}^\top \mathbf{A} \mathbf{v}$  is always an *eigenvalue* of  $\mathbf{A}$ , and the vector  $\mathbf{v}$  that achieves this maximum value is always an *eigenvector* corresponding to that eigenvalue. Below we review the definition of eigenvalues and eigenvectors, and then give the statement of the Spectral Theorem.

(1.34) **DEFINITION**

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a square matrix. For  $\mathbf{v} \in \mathbb{R}^n$  and  $\lambda \in \mathbb{R}$ , we say that  $\mathbf{v}$  is an *eigenvector* of  $\mathbf{A}$  with *eigenvalue*  $\lambda$  if

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}.$$

(1.35) **THEOREM (Spectral Theorem)**

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a symmetric matrix. Then, there is an orthonormal basis  $\{\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(n)}\}$  for  $\mathbb{R}^n$  that consists entirely of eigenvectors of  $\mathbf{A}$ , with associated eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_n$ , respectively. Assuming without loss of generality that  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ , we have

$$(\mathbf{v}^{(1)})^\top \mathbf{A} \mathbf{v}^{(1)} = \lambda_1 = \max_{\mathbf{v}} \mathbf{v}^\top \mathbf{A} \mathbf{v} \text{ subject to } \mathbf{v} \neq \mathbf{0},$$

and for  $i > 1$  we have

$$(\mathbf{v}^{(i)})^\top \mathbf{A} \mathbf{v}^{(i)} = \lambda_i = \max_{\mathbf{v}} \mathbf{v}^\top \mathbf{A} \mathbf{v} \text{ subject to } \mathbf{v} \neq \mathbf{0} \text{ and } \forall j < i [\mathbf{v}^\top \mathbf{v}^{(j)} = 0].$$

Applying the Spectral Theorem to  $\text{cov}(\mathbf{X})$ , we discover that the direction of greatest variance—that is, the unit vector that maximizes  $\text{cov}(\mathbf{X}\mathbf{v}) = \mathbf{v}^\top \text{cov}(\mathbf{X})\mathbf{v}$ —is the unit-length eigenvector for  $\mathbf{X}$  with the greatest corresponding eigenvalue. Therefore, our principal components  $\mathbb{B}$  are simply an orthonormal set of eigenvectors of  $\mathbf{A}$ , where the  $i$ th principal component  $\psi^{(i)}$  is the eigenvector with the  $i$ th greatest eigenvalue.

**Finding the Eigenvectors.** The final piece of the PCA algorithm is to find a way to compute the eigenvectors of the covariance matrix. Typically, this is achieved using a technique called *singular value decomposition* (SVD).

(1.36) **THEOREM** (Singular Value Decomposition)

Let  $A \in \mathbb{R}^{m \times n}$  be a matrix, not necessarily square. Then, there exist matrices  $U \in \mathbb{R}^{m \times m}$ ,  $S \in \mathbb{R}^{m \times n}$ , and  $V \in \mathbb{R}^{n \times n}$  such that

$$A = USV^T$$

and the following hold.

- The columns of  $U$  and  $V$  are pairwise orthogonal and of unit length.
- $S$  is diagonal:  $S_{i,j} = 0$  whenever  $i \neq j$ .
- The diagonal entries of  $S$  are in decreasing order: if  $i < j \leq \min(m, n)$ , then  $S_{i,i} \geq S_{j,j}$ .
- Each column of  $V$  is an eigenvector of  $A^T A$ . The eigenvalue corresponding to the  $i$ th column of  $V$  if  $i \leq \min(m, n)$  and 0 otherwise.
- Each column of  $U$  is an eigenvector of  $AA^T$ . The eigenvalue corresponding to the  $i$ th column of  $U$  is  $S_{i,i}^2$  if  $i \leq \min(m, n)$  and 0 otherwise.

There are already algorithms for computing the SVD of a matrix. For PCA, we apply SVD to the matrix  $A$  defined as follows:

$$A = \frac{X - \bar{X}}{\sqrt{m}}.$$

Observe that  $\text{cov}(X) = A^T A$ . If  $A = USV^T$  is the SVD of  $A$ , then the columns of  $V$  contain pairwise-orthogonal unit eigenvectors of  $\text{cov}(X)$ , sorted in decreasing order of their associated eigenvalues. This means that the  $i$ th column of  $V$  is in fact the  $i$ th principal component of  $X$ , which makes  $V^T$  the change of basis matrix for  $\mathbb{B}$ .

**PCA Procedure.** To summarize, the full procedure for embedding vectors into a lower-dimensional space using PCA is as follows.

**ALGORITHM: DIMENSIONALITY REDUCTION VIA PRINCIPAL COMPONENTS ANALYSIS**

**Inputs:**

- $X \in \mathbb{R}^{m \times n}$ , a set of feature vectors arranged into a matrix of row vectors
- $k$ , the desired dimensionality of the embedding vectors

**Output:**

- $X' \in \mathbb{R}^{m \times k}$ , an embedding of  $X$  into  $\mathbb{R}^k$

**Procedure:**

- Compute the mean of  $X$ :

$$\bar{X} \leftarrow \frac{X^T \mathbf{1}}{m}.$$

- Compute the “square root” of the covariance matrix of  $\mathbf{X}$ :

$$\mathbf{A} \leftarrow \frac{(\mathbf{X} - \bar{\mathbf{X}})}{\sqrt{m}}.$$

- Compute the SVD of  $\mathbf{A}$ : let  $\mathbf{U}$ ,  $\mathbf{S}$ , and  $\mathbf{V}$  be such that

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$$

is the SVD of  $\mathbf{A}$ .

- Truncate  $\mathbf{V}$  to include only the first  $k$  principal components of  $\mathbf{X}$ :

$$\mathbf{V} \leftarrow \mathbf{V}_{:,k}.$$

- Project  $\mathbf{X}$  onto its first  $k$  principal components using  $\mathbf{V}^\top$  as a change of basis matrix:

$$\mathbf{X}' \leftarrow \mathbf{X}\mathbf{V}.$$

- Return  $\mathbf{X}'$ .

## 1.5.2 t-Distributed Stochastic Neighbors Embedding

The  $t$ -SNE algorithm embeds a finite set  $\mathbb{D} \subseteq \mathbb{R}^n$  of feature vectors into a lower-dimensional space  $\mathbb{R}^k$  in a way that is isometric with respect to two similarity metrics,  $d_1 : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  and  $d_2 : \mathbb{R}^k \times \mathbb{R}^k \rightarrow \mathbb{R}$ . These two similarity metrics are based on the probability that two randomly chosen points are “neighbors.” In the high-dimensional space,  $d_1$  is based on a normal distribution, while in the low-dimensional space,  $d_2$  is based on the Student’s  $t$  distribution with 1 degree of freedom. The embedding is computed by viewing  $d_1$  and  $d_2$  as probability distributions and maximizing their similarity.

**High-Dimensional Similarity.** Fix a point  $\mathbf{x} \in \mathbb{D}$ . Imagine that we want to randomly choose a “neighbor”  $\mathbf{y}$  of  $\mathbf{x}$ , in a way such that the Euclidean distance  $\|\mathbf{y} - \mathbf{x}\|$  between  $\mathbf{y}$  and  $\mathbf{x}$  is normally distributed with standard deviation  $\sigma_{\mathbf{x}}$ . The probability density function for  $\|\mathbf{y} - \mathbf{x}\|$  is given by

$$p_{\mathbf{x}}(\mathbf{y}) = \frac{1}{\sigma_{\mathbf{x}}\sqrt{2\pi}} e^{-\|\mathbf{y}-\mathbf{x}\|^2/(2\sigma_{\mathbf{x}}^2)}.^{14}$$

Observe that points closer to  $\mathbf{x}$ , which are more likely to be considered “neighbors” of  $\mathbf{x}$ , have a higher density than points that are farther away from  $\mathbf{x}$ .

Now, imagine that we are choosing a random point of  $\mathbb{D}$ , not including  $\mathbf{x}$ , to be a “neighbor” of  $\mathbf{x}$ . We want this random choice to be distributed according to our normal distribution. However,

<sup>14</sup>Technically, this implies that  $\|\mathbf{y} - \mathbf{x}\|$  can be negative; we won’t think about this too much.



the normal distribution is a continuous distribution, while sampling from  $\mathbb{D}$  requires a discrete distribution. In order to turn our normal distribution into a discrete distribution, we convert  $p_x$  into a probability mass function by normalizing over  $\mathbb{D} \setminus \{\mathbf{x}\}$ :

$$\mathbb{P}[\mathbf{y} \in \mathbb{D} \setminus \{\mathbf{x} \text{ is chosen as a "neighbor" of } \mathbf{x}\}] = p(\mathbf{y}, \mathbf{x}) = \frac{p_x(\mathbf{y})}{\sum_{\mathbf{z} \in \mathbb{D} \setminus \{\mathbf{x}\}} p_x(\mathbf{z})}.$$

Our similarity metric on  $\mathbb{R}^n$  is defines the similarity between  $\mathbf{x}$  and  $\mathbf{y}$  as the average of  $p(\mathbf{x}, \mathbf{y})$  and  $p(\mathbf{y}, \mathbf{x})$ , scaled by the cardinality of  $\mathbb{D}$ :

$$d_1(\mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y}) + p(\mathbf{y}, \mathbf{x})}{2|\mathbb{D}|}.$$

Assuming that  $d(\mathbf{x}, \mathbf{y}) = 0$  when  $\mathbf{x} = \mathbf{y}$ , observe that  $d_1$  is a probability mass function over  $\mathbb{D} \times \mathbb{D}$ :

$$\sum_{\mathbf{x}, \mathbf{y} \in \mathbb{D}} d_1(\mathbf{x}, \mathbf{y}) = \frac{1}{2|\mathbb{D}|} \left( \left( \sum_{\mathbf{x}, \mathbf{y} \in \mathbb{D}} p(\mathbf{x}, \mathbf{y}) \right) + \left( \sum_{\mathbf{x}, \mathbf{y} \in \mathbb{D}} p(\mathbf{y}, \mathbf{x}) \right) \right) = \frac{1}{2|\mathbb{D}|} (|\mathbb{D}| + |\mathbb{D}|) = 1.$$

**Low-Dimensional Similarity.** To define  $d_2$ , we do essentially the same thing, except we use a Student's  $t$ -distribution with 1 degree of freedom centered around  $\mathbf{x}$ . The probability density function for this distribution is given by

$$q_x(\mathbf{y}) = \frac{1}{\pi(1 + \|\mathbf{y} - \mathbf{x}\|^2)}.$$

Unlike  $p_x$ ,  $q_x$  is “symmetric” in the sense that  $q_x(\mathbf{y}) = q_y(\mathbf{x})$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ . Therefore, rather than viewing  $q_x(\mathbf{y})$  as the process of choosing a neighbor  $\mathbf{y}$  for a pre-determined point  $\mathbf{x}$ , we can view  $q_x(\mathbf{y})$  as simultaneously choosing two points,  $\mathbf{x}$  and  $\mathbf{y}$ , to be neighbors. Thus, given a set of embeddings  $\mathbb{E}$ , we define a probability mass function over  $\mathbb{E}$  by

$$\mathbb{P}[\mathbf{x} \text{ and } \mathbf{y} \text{ are "neighbors"}] = q(\mathbf{x}, \mathbf{y}) = \frac{q_x(\mathbf{y})}{\sum_{\mathbf{u}, \mathbf{v} \in \mathbb{E}, \mathbf{u} \neq \mathbf{v}} q_u(\mathbf{v})},$$

and we simply set  $d_2(\mathbf{x}, \mathbf{y}) = q(\mathbf{x}, \mathbf{y})$  when  $\mathbf{x} \neq \mathbf{y}$  and  $d_2(\mathbf{x}, \mathbf{y}) = 0$  otherwise. Note that  $d_2$  is also a valid probability mass function.

**Finding an Optimal Embedding.** Let  $f^* : \mathbb{R}^n \rightarrow \mathbb{R}^k$  be our  $t$ -SNE embedding. We would like  $f^*$  to embed  $\mathbb{D}$  as isometrically as possible with respect to  $d_1$  and  $d_2$ . The degree of isometricity is measured by the *Kullback–Leibler divergence* between the probability distribution  $d_1$  and the probability distribution defined by  $d_2$ .

(1.37) **DEFINITION**

Let  $p, q : \mathbb{A} \rightarrow [0, 1]$  be probability mass functions over a set  $\mathbb{A}$ . The *Kullback–Leibler divergence* of  $p$  from  $q$ , or *KL-divergence* for short, is defined as

$$\text{KL}(p \| q) = \sum_{a \in \mathbb{A}} p(a) \ln \left( \frac{p(a)}{q(a)} \right).$$

Note that the definition of  $d_2$  depends on how  $\mathbb{D}$  is embedded into  $\mathbb{R}^k$ . Therefore, different embeddings  $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$  will result in different definitions of  $d_2$ . We define the *error* of an embedding  $f$  to be the KL-divergence between  $d_1$  and the version of  $d_2$  given by  $f$ :

$$E(f) = \sum_{\mathbf{x}, \mathbf{y} \in \mathbb{D}} d_1(\mathbf{x}, \mathbf{y}) \ln \left( \frac{d_1(\mathbf{x}, \mathbf{y})}{d_2(f(\mathbf{x}), f(\mathbf{y}))} \right).$$

We then compute our  $t$ -SNE embedding by minimizing  $E(f)$ :

$$f^* = \operatorname{argmin}_f E(f) = \operatorname{argmin}_f \sum_{\mathbf{x}, \mathbf{y} \in \mathbb{D}} d_1(\mathbf{x}, \mathbf{y}) \ln \left( \frac{d_1(\mathbf{x}, \mathbf{y})}{d_2(f(\mathbf{x}), f(\mathbf{y}))} \right).$$

# Bibliography

- Leonard Bloomfield. 1933. *Language*, first edition. Henry Holt and Company, New York, NY.
- J. R. Firth. 1957. A Synopsis of Linguistic Theory, 1930–1955. In *Studies in Linguistic Analysis*, first edition, Special Volume of the Philological Society, pages 1–31. Blackwell Publishers, Oxford, United Kingdom.
- Yoav Goldberg and Omer Levy. 2014. Word2vec Explained: Deriving Mikolov et al.’s Negative-Sampling Word-Embedding Method. *Computing Research Repository*, arXiv:1402.3722.
- Thomas Kollar, Danielle Berry, Lauren Stuart, Karolina Owczarzak, Tagyoung Chung, Lambert Mathias, Michael Kayser, Bradford Snow, and Spyros Matsoukas. 2018. The Alexa Meaning Representation Language. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, volume 3 (industry papers), pages 177–184, New Orleans, LA. Association for Computational Linguistics.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient Estimation of Word Representations in Vector Space. In *ICLR 2013 Workshop Proceedings*, Scottsdale, AZ. OpenReview.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed Representations of Words and Phrases and Their Compositionality. In *Advances in Neural Information Processing Systems 26*, pages 3111–3119, Montreal, Canada. Curran Associates, Inc.