

## Preface

My interest in text began in the late 1990s during my teenage years, building dynamic websites using Perl and HTML. This early experience with coding and organizing text into structured formats sparked my fascination with how text could be processed and transformed. Over the years, I advanced to building web scrapers and text aggregators, developing systems to extract structured data from webpages. The challenge of processing and understanding text led me to explore more complex applications, including designing chatbots that could understand and address user needs.

The challenge of extracting meaning from words intrigued me. The complexity of the task only fueled my determination to “crack” it, using every tool at my disposal—ranging from regular expressions and scripting languages to text classifiers and named entity recognition models.

The rise of large pretrained language models (LLMs) transformed everything. For the first time in history, computers could converse with us fluently and follow verbal instructions with remarkable precision. However, like any tool, their immense power comes with limitations. Some are easy to spot, but others are more subtle, requiring deep expertise to handle properly. Attempting to build a skyscraper without fully understanding your tools will only result in a pile of concrete and steel. The same holds true for language models. Approaching large-scale text processing tasks or creating reliable products for paying users requires precision and knowledge—guesswork simply isn’t an option.

To make this book engaging and to deepen the reader’s understanding, I decided to discuss language modeling as a whole, including approaches that are often overlooked in modern literature. While LLMs dominate the spotlight, older approaches like count-based methods and recurrent neural networks (RNNs) remain highly effective, especially when efficiency and cost are priorities. By revisiting these foundational methods, my goal is to highlight the field’s evolution and equip readers with a broader set of tools for solving real-world problems.

I wrote this book for those who, like me, are captivated by the challenge of understanding language through machines. Language models are, at their core, just mathematical functions. However, their true potential isn’t fully appreciated in theory—you need to implement them to see their power and how their abilities grow as they scale. This is why I decided to make this book hands-on, unlike my previous two.

As language modeling continues to evolve, so do the challenges it presents. By the end of this book, I hope you’ll not only have enjoyed the process but also feel confident in applying LLMs and other techniques effectively, with a solid grasp of both their “magic” and their limitations.

If this is your first time exploring language models, I envy you a little—it’s truly magical to discover how machines learn to understand the world through natural language.

I hope you enjoy reading this book as much as I enjoyed writing it!

## About the Book

This book provides a solid foundation in language modeling fundamentals. Instead of covering every recent development—which can quickly become outdated—it emphasizes core concepts that remain useful across a variety of contexts. The emphasis on fundamentals helps readers build lasting knowledge that can be applied even as specific techniques evolve.

The content moves from early count-based models to recurrent neural networks and transformers. Its hands-on approach, supported by mathematics and illustrations, helps readers build confidence in their tools and prepares them to work with any modern language model.

## Who This Book Is For

This book serves software developers, data scientists, machine learning engineers, and anyone curious about language models. Whether you’re integrating existing models into applications or training your own, you’ll find practical guidance alongside theoretical foundations.

Given its hundred-page format, the book makes certain assumptions about readers. You should have programming experience, as all hands-on examples use Python.

While familiarity with PyTorch and tensors—PyTorch’s fundamental data types—is beneficial, it’s not mandatory. If you’re new to these tools, the book’s wiki provides a concise introduction with examples and resource links for further learning. This wiki format ensures content remains current and addresses reader questions beyond publication.

College-level math knowledge helps, but you needn’t remember every detail or have machine learning experience. The book introduces concepts systematically, beginning with notations, definitions, and fundamental vector and matrix operations. From there, it progresses through simple neural networks to more advanced topics. Mathematical concepts are presented intuitively, with clear diagrams and examples that facilitate understanding.

## What This Book Is Not

This book is focused on understanding and implementing language models. It will *not* cover:

- **Large-scale training:** This book won’t teach you how to train massive models on distributed systems or how to manage training infrastructure.
- **Production deployment:** Topics like model serving, API development, scaling for high traffic, monitoring, and cost optimization are not covered. The code examples focus on understanding the concepts rather than production readiness.
- **Enterprise applications:** This book won’t guide you through building commercial LLM applications, handling user data, or integrating with existing systems.

If you’re interested in learning the mathematical foundations of language models, understanding how they work, implementing core components yourself, or learning to work effectively with LLMs,

this book is for you. But if you’re primarily looking to deploy models in production or build scalable applications, you may want to supplement this book with other resources.

## Book Structure

The book is divided into six chapters, progressing from fundamentals to advanced topics:

- **Chapter 1** covers machine learning basics, including key concepts like AI, models, neural networks, and gradient descent. Even if you’re familiar with these topics, the chapter provides important foundations for understanding language models.
- **Chapter 2** introduces language modeling fundamentals, exploring text representation methods like bag of words and word embeddings, as well as count-based language models and evaluation techniques.
- **Chapter 3** focuses on recurrent neural networks, covering their implementation, training, and application as language models.
- **Chapter 4** provides a detailed exploration of the Transformer architecture, including key components like self-attention, position embeddings, and practical implementation.
- **Chapter 5** examines large language models (LLMs), discussing why scale matters, finetuning techniques, practical applications, and important considerations around hallucinations, copyright, and ethics.
- **Chapter 6** concludes with further reading on advanced topics like mixture of experts, model compression, preference-based alignment, and vision language models, providing direction for continued learning.

Most chapters contain working code examples you can run and modify. While only essential code appears in the book, complete code is available as Jupyter notebooks on the book’s website, with notebooks referenced in relevant sections. All code in notebooks remains compatible with the latest stable versions of Python, PyTorch, and other libraries.

The notebooks run on Google Colab, which at the time of writing offers free access to computing resources including GPUs and TPUs. These resources, though, aren’t guaranteed and have usage limits that may vary. Some examples might require extended GPU access, potentially involving wait times for availability. If the free tier proves limiting, Colab’s pay-as-you-go option lets you purchase compute credits for reliable GPU access. While these credits are relatively affordable by North American standards, costs may be significant depending on your location.

For those familiar with the Linux command line, GPU cloud services provide another option through pay-per-time virtual machines with one or more GPUs. The book’s wiki maintains current information on free and paid notebook or GPU rental services.

**Verbatim** terms and blocks indicate code, code fragments, or code execution outputs. **Bold** terms link to the book’s term index, and occasionally highlight algorithm steps.

## Should You Buy This Book?

Like my previous two books, this one is distributed on the *read first, buy later* principle. I firmly believe that paying for content before consuming it means buying a pig in a poke. At a dealership, you can see and try a car. In a department store, you can try on clothes. Similarly, you should be able to read a book before paying for it.

The *read first, buy later* principle means you can freely download the book, read it, and share it with friends and colleagues. If you find the book helpful or useful in your work, business, or studies—or if you simply enjoy reading it—then buy it.

Now grab your tea or coffee, and let's begin!

## Acknowledgements

The high quality of this book would be impossible without volunteering editors. I especially thank Erman Sert, Viet Hoang Tran Duong, Alex Sherstinsky, Kelvin Sundli, and Mladen Korunoski for their systematic contributions.

I am also grateful to Alireza Bayat Makou, Taras Shalaiko, Domenico Siciliani, Preethi Raju, Sri Kumar Sundareshwar, Mathieu Nayrolles, Abhijit Kumar, Giorgio Mantovani, Abhinav Jain, Steven Finkelstein, Ryan Gaughan, Ankita Guha, Harmanan Kohli, Daniel Gross, Kea Kohv, Marcus Oliveira, Tracey Mercier, Prabin Kumar Nayak, Saptarshi Datta, Gurgen R. Hayrapetyan, Sina Abdidizaji, Federico Raimondi Cominesi, and Manoj Pillai for their help.

# Chapter 1. Machine Learning Basics

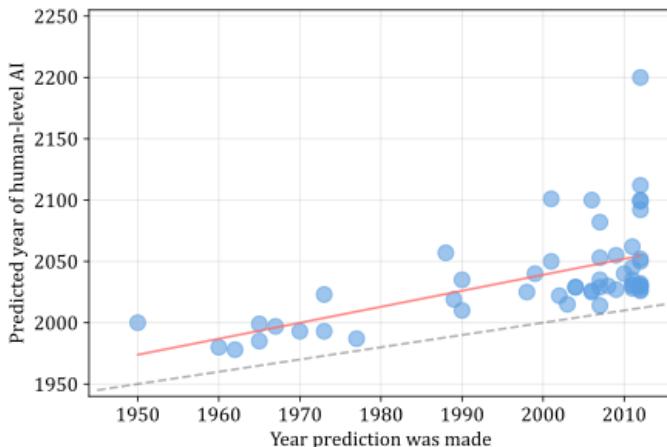
This chapter introduces the fundamental concepts of machine learning. Starting with the evolution of artificial intelligence and machine learning, it defines a model and presents the four-step machine learning process. The chapter then covers essential mathematical foundations, including vectors and matrices, before examining neural networks. It concludes with key optimization techniques, focusing on gradient descent and automatic differentiation.

## 1.1. AI and Machine Learning

The term “artificial intelligence” (AI) was first introduced in 1955 during a workshop led by John McCarthy, focusing on exploring how machines could use language, form concepts, solve problems like humans, and improve over time. Building on these ideas, Joseph Weizenbaum developed the first chatbot, ELIZA, in 1966. ELIZA simulated conversations by detecting patterns in user input and replying with preprogrammed responses, giving the impression of understanding.

In AI’s early years, researchers were overly optimistic about achieving human-level intelligence. In 1965, Herbert Simon, a Turing Award recipient, predicted that “machines will be capable, within twenty years, of doing any work a man can do.” However, progress was slower than expected, leading to periods of reduced funding and interest, known as “AI winters.”

Interestingly, since the 1950s, experts consistently predicted that human-level AI will be achieved in about 25 years:



Two major AI winters occurred in 1974–1980 and 1987–2000. These periods were marked by notable setbacks, including the failure of machine translation in 1966, poor outcomes from DARPA’s Speech Understanding Research program at Carnegie Mellon (1971–1975), and reduced AI funding in the UK after the 1973 Lighthill report. In the 1990s, many expert systems—computer programs that simulated human decision-making using predefined rules and domain-specific logic—were abandoned due to high costs and limited success.

During the first AI winter, even the term “AI” became somewhat taboo. Many researchers rebranded their work as “informatics,” “knowledge-based systems,” or “pattern recognition” to avoid association with AI’s perceived failures.

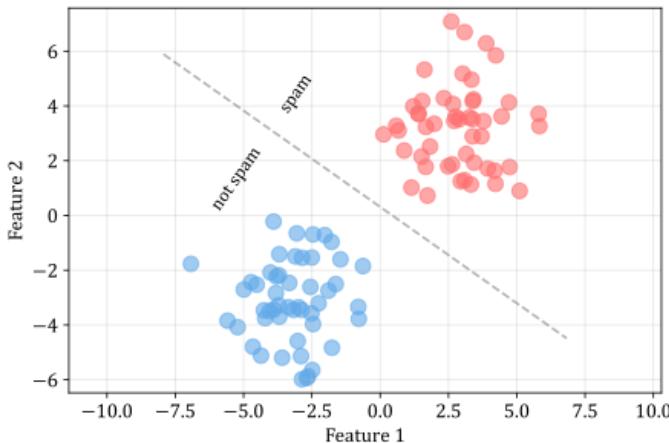
Enthusiasm for AI has grown steadily since the early 1990s. Interest surged around 2012, particularly in machine learning, driven by advances in computational power, access to large datasets, and improvements in neural network algorithms and frameworks. These developments led to increased funding and a significant AI boom.

Although the focus of **artificial intelligence** research has evolved, the core goal remains the same: to create methods that enable machines to solve problems previously considered solvable only by humans. This is how the term will be used throughout this book.

The term “machine learning” was introduced in 1959 by Arthur Samuel. In his paper, “Some Studies in Machine Learning Using the Game of Checkers,” he described it as “programming computers to learn from experience.”

Early AI researchers primarily focused on symbolic methods and rule-based systems—an approach later dubbed **good old-fashioned AI (GOFAI)**—but over time, the field increasingly embraced machine learning approaches, with neural networks emerging as a particularly powerful technique.

Neural networks, inspired by the brain, aimed to learn patterns directly from examples. One foundational model, the perceptron, was introduced by Frank Rosenblatt in 1958. It became a key step toward later advancements. The perceptron defines a decision boundary, a line that separates examples of two classes (e.g., spam and not spam):



Decision trees and random forests represent important evolutionary steps in machine learning. **Decision trees**, introduced in the 1960s and later advanced by Ross Quinlan’s ID3 algorithm in 1986, split data into subsets through a tree-like structure. Each node represents a question about the data, each branch is an answer, and each leaf provides a prediction. While these models are easy

to understand, they can struggle with **overfitting**, where they adapt too closely to training data, reducing their ability to perform well on new, unseen data.

To address this limitation, Leo Breiman introduced the random forest algorithm in 2001. A **random forest** builds multiple decision trees using random subsets of data and combines their outputs. This approach improves predictive accuracy and reduces overfitting. Random forests remain widely used for their reliability and performance.

**Support vector machines (SVMs)**, introduced in the 1990s by Vladimir Vapnik and his colleagues, were another significant step forward. SVMs identify the optimal hyperplane that separates data points of different classes with the widest margin. The introduction of **kernel methods** allowed SVMs to manage complex, non-linear patterns by mapping data into higher-dimensional spaces, making it easier to find a suitable separating hyperplane. These advances made SVMs central to machine learning research.

Today, **machine learning** is a subfield of AI focused on creating algorithms that learn from collections of examples. These examples can come from nature, be designed by humans, or be generated by other algorithms. The process involves gathering a dataset and building a model from it, which is then used to solve the problem.

I will use “learning” and “machine learning” interchangeably to save keystrokes.

## 1.2. Model

A **model** is typically represented by a mathematical equation:

$$y = f(x)$$

Here,  $x$  is the input,  $y$  is the output, and  $f$  represents a function of  $x$ . A **function** is a named rule that describes how one set of values is related to another. Formally, a function  $f$  maps inputs from the **domain** to outputs in the **codomain**, ensuring each input has exactly one output. The function uses a specific rule or formula to transform the input into the output.

In machine learning, the goal is to compile a **dataset** of **examples** and use them to build  $f$ , so when  $f$  is applied to a new, unseen  $x$ , it produces a  $y$  that gives meaningful insight into  $x$ .

To predict a house’s price based on its area, the dataset might include (area, price) pairs such as  $\{(150,200), (200,600), \dots\}$ . Here, the area is measured in  $\text{m}^2$ , and the price is in thousands.

Curly brackets denote a set. A set containing  $N$  elements, ranging from  $x_1$  to  $x_N$ , is expressed as  $\{x_i\}_{i=1}^N$ .

Imagine we own a house with an area of  $250 \text{ m}^2$  (about 2691 square feet). To derive a function  $f$  that provides a reasonable price for this house, testing every possible function is infeasible. Instead, we select a specific *structure* for  $f$  and focus on functions that match this structure.

Let's define the structure for  $f$  as:

$$f(x) \stackrel{\text{def}}{=} wx + b, \quad (1.1)$$

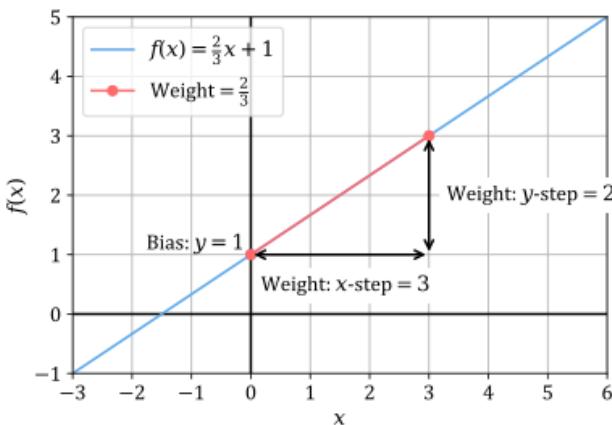
which describes a **linear function** of  $x$ . The formula  $wx + b$  is a **linear transformation** of  $x$ .

The notation  $\stackrel{\text{def}}{=}$  means "equals by definition" or "is defined as."

For linear functions, determining  $f$  requires only two values:  $w$  and  $b$ . These are called the **parameters** or **weights** of the model.

In other texts,  $w$  might be referred to as the **slope**, **coefficient**, or **weight term**. Similarly,  $b$  may be called the **intercept**, **constant term**, or **bias**. In this book, we'll stick to "weight" for  $w$  and "bias" for  $b$ , as these terms are widely used in machine learning. When the meaning is clear, "parameters" and "weights" will be used interchangeably.

For instance, when  $w = \frac{2}{3}$  and  $b = 1$ , the linear function is shown below:



Here, the bias shifts the graph vertically, so the line crosses the  $y$ -axis at  $y = 1$ . The weight determines the slope, meaning the line rises by 2 units for every 3 units it moves to the right.

Mathematically, the function  $f(x) = wx + b$  is an **affine transformation**, rather than a linear one, since a true linear transformation requires  $b = 0$ . In machine learning, however, we often call such models "linear" as long as the parameters appear linearly in the equation. This means  $w$  and  $b$  are only multiplied by the inputs or constants and added—they don't multiply each other, get raised to powers, or appear inside functions like  $\sin(w)$  or  $e^b$ .

Even with a simple model like  $f(x) = wx + b$ , the parameters  $w$  and  $b$  can take infinitely many values. To find the optimal ones, we need an optimality criterion. A reasonable choice is to minimize the average error when predicting house prices based on area. In this case, we aim for  $f(x) = wx + b$  to make predictions as close as possible to the actual prices.

Let our dataset be  $\{(x_i, y_i)\}_{i=1}^N$ , where  $N$  is the size of the dataset and  $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$  are individual examples. In machine learning,  $x_i$  is called the **input**, and  $y_i$  is the **target**. When every example includes both an input and a target, the learning process is known as **supervised**. This book's focus is supervised machine learning.

Other machine learning types include **unsupervised learning**, where models learn patterns from inputs alone, and **reinforcement learning**, where models learn by interacting with environments and receiving rewards or penalties for their actions.

When  $f(x)$  is applied to  $x_i$ , it generates a predicted value  $\tilde{y}_i$ . We can define the error  $\text{err}(\tilde{y}_i, y_i)$  for a given example  $(x_i, y_i)$  as:

$$\text{err}(\tilde{y}_i, y_i) \stackrel{\text{def}}{=} (\tilde{y}_i - y_i)^2 \quad (1.2)$$

This expression, called the **squared error**, equals 0 when  $\tilde{y}_i = y_i$ . This makes sense: there's no error if the predicted price matches the actual price. The further  $\tilde{y}_i$  deviates from  $y_i$ , the larger the error becomes. Squaring ensures the error is always positive, whether the prediction overshoots or undershoots.

We define  $w^*$  and  $b^*$  as the optimal parameter values for  $f$ , which minimize the average price prediction error for the dataset using the following expression:

$$\frac{\text{err}(\tilde{y}_1, y_1) + \text{err}(\tilde{y}_2, y_2) + \dots + \text{err}(\tilde{y}_N, y_N)}{N}$$

Let's rewrite the above expression by expanding each  $\text{err}(\cdot)$ :

$$\frac{(\tilde{y}_1 - y_1)^2 + (\tilde{y}_2 - y_2)^2 + \dots + (\tilde{y}_N - y_N)^2}{N}$$

Let's assign the name  $J(w, b)$  to our expression, turning it into a function:

$$J(w, b) \stackrel{\text{def}}{=} \frac{(wx_1 + b - y_1)^2 + (wx_2 + b - y_2)^2 + \dots + (wx_N + b - y_N)^2}{N} \quad (1.3)$$

In the equation defining  $J(w, b)$ , which represents the average prediction error, the values of  $x_i$  and  $y_i$  for each  $i$  from 1 to  $N$  are known since they come from the dataset. The unknowns are  $w$  and  $b$ . To determine the optimal  $w^*$  and  $b^*$ , we need to minimize  $J(w, b)$ . As this function is quadratic in two variables, calculus guarantees it has a single minimum.

The expression Equation 1.3 is referred to as the **loss function** in the machine learning problem of **linear regression**. In this case, the loss function is the **mean squared error** or **MSE**.

To find the optimum (minimum or maximum) of a function, we calculate its **first derivative**. When we reach the optimum, the first derivative equals zero. For functions of two or more variables, like the loss function  $J(w, b)$ , we compute **partial derivatives** with respect to each variable. We denote these as  $\frac{\partial J}{\partial w}$  for  $w$  and  $\frac{\partial J}{\partial b}$  for  $b$ .

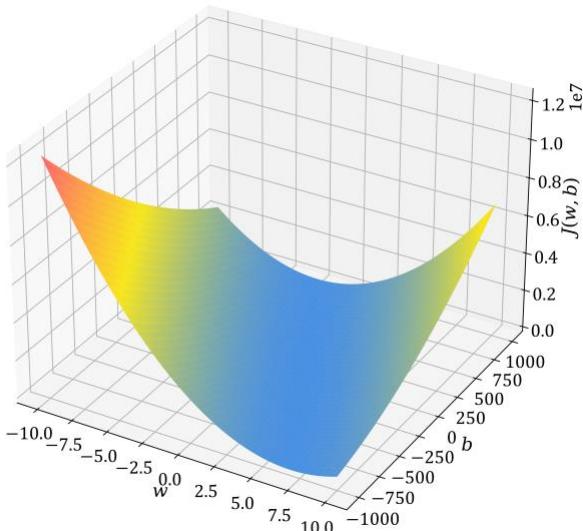
To determine  $w^*$  and  $b^*$ , we solve the following system of two equations:

$$\begin{cases} \frac{\partial J}{\partial w} = 0 \\ \frac{\partial J}{\partial b} = 0 \end{cases}$$

Fortunately, the mean squared error function's structure and the model's linearity allow us to solve this system of equations analytically. To illustrate, consider a dataset with three examples:  $(x_1, y_1) = (150, 200)$ ,  $(x_2, y_2) = (200, 600)$ , and  $(x_3, y_3) = (260, 500)$ . For this dataset, the loss function is:

$$J(w, b) \stackrel{\text{def}}{=} \frac{(150w + b - 200)^2 + (200w + b - 600)^2 + (260w + b - 500)^2}{3}$$

Let's plot it:



Navigate to the book's wiki, from the file [thelmbook.com/py/1.1](http://thelmbook.com/py/1.1) retrieve the code used to generate the above plot, run the code, and rotate the graph in 3D to better observe the minimum.

Now we need to derive the expressions for  $\frac{\partial J}{\partial w}$  and  $\frac{\partial J}{\partial b}$ . Notice that  $J(w, b)$  is a composition of the following functions:

- Functions  $d_1 \stackrel{\text{def}}{=} 150w + b - 200$ ,  $d_2 \stackrel{\text{def}}{=} 200w + b - 600$ ,  $d_3 \stackrel{\text{def}}{=} 260w + b - 500$  are linear functions of  $w$  and  $b$ ;
- Functions  $\text{err}_1 \stackrel{\text{def}}{=} d_1^2$ ,  $\text{err}_2 \stackrel{\text{def}}{=} d_2^2$ ,  $\text{err}_3 \stackrel{\text{def}}{=} d_3^2$  are quadratic functions of  $d_1$ ,  $d_2$ , and  $d_3$ ;
- Function  $J \stackrel{\text{def}}{=} \frac{1}{3}(\text{err}_1 + \text{err}_2 + \text{err}_3)$  is a linear function of  $\text{err}_1$ ,  $\text{err}_2$ , and  $\text{err}_3$ .

**A composition of functions** means the output of one function becomes the input to another. For example, with two functions  $f$  and  $g$ , you first apply  $g$  to  $x$ , then apply  $f$  to the result. This is written as  $f(g(x))$ , which means you calculate  $g(x)$  first and then use that result as the input for  $f$ .

In our loss function  $J(w, b)$ , the process starts by computing the linear functions for  $d_1$ ,  $d_2$ , and  $d_3$  using the current values of  $w$  and  $b$ . These outputs are then passed into the quadratic functions  $\text{err}_1$ ,  $\text{err}_2$ , and  $\text{err}_3$ . The final step is averaging these results to compute  $J$ .

Using the sum rule and the constant multiple rule of differentiation,  $\frac{\partial J}{\partial w}$  is given by:

$$\frac{\partial J}{\partial w} = \frac{1}{3} \left( \frac{\partial \text{err}_1}{\partial w} + \frac{\partial \text{err}_2}{\partial w} + \frac{\partial \text{err}_3}{\partial w} \right),$$

where  $\frac{\partial \text{err}_1}{\partial w}$ ,  $\frac{\partial \text{err}_2}{\partial w}$ , and  $\frac{\partial \text{err}_3}{\partial w}$  are the partial derivatives of  $\text{err}_1$ ,  $\text{err}_2$ , and  $\text{err}_3$  with respect to  $w$ .

The **sum rule** of differentiation states that the derivative of the sum of two functions equals the sum of their derivatives:  $\frac{\partial}{\partial x}[f(x) + g(x)] = \frac{\partial}{\partial x}f(x) + \frac{\partial}{\partial x}g(x)$ .

The **constant multiple rule** of differentiation states that the derivative of a constant multiplied by a function equals the constant times the derivative of the function:  $\frac{\partial}{\partial x}[c \cdot f(x)] = c \cdot \frac{\partial}{\partial x}f(x)$ .

By applying the chain rule of differentiation, the partial derivatives of  $\text{err}_1$ ,  $\text{err}_2$ , and  $\text{err}_3$  with respect to  $w$  are:

$$\begin{aligned} \frac{\partial \text{err}_1}{\partial w} &= \frac{\partial \text{err}_1}{\partial d_1} \cdot \frac{\partial d_1}{\partial w}, && \text{partial derivative of } d_1 \\ &\quad \text{with respect to } w \\ \frac{\partial \text{err}_2}{\partial w} &= \frac{\partial \text{err}_2}{\partial d_2} \cdot \frac{\partial d_2}{\partial w}, && \text{multiplied by} \\ \frac{\partial \text{err}_3}{\partial w} &= \frac{\partial \text{err}_3}{\partial d_3} \cdot \frac{\partial d_3}{\partial w} \end{aligned}$$

The **chain rule** of differentiation states that the derivative of a **composite function**  $f(g(x))$ , written as  $\frac{\partial}{\partial x} [f(g(x))]$ , is the product of the derivative of  $f$  with respect to  $g$  and the derivative of  $g$  with respect to  $x$ , or:  $\frac{\partial}{\partial x} [f(g(x))] = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x}$ .

Then,

$$\begin{aligned} \frac{\partial \text{err}_1}{\partial d_1} &= \frac{\partial \text{err}_1}{\partial w} = 2d_1 \cdot 150 = 300 \cdot (150w + b - 200), \\ \frac{\partial \text{err}_2}{\partial d_2} &= \frac{\partial \text{err}_2}{\partial w} = 2d_2 \cdot 200 = 400 \cdot (200w + b - 600), \\ \frac{\partial \text{err}_3}{\partial d_3} &= \frac{\partial \text{err}_3}{\partial w} = 2d_3 \cdot 260 = 520 \cdot (260w + b - 500) \end{aligned}$$

Therefore,

$$\begin{aligned} \frac{\partial J}{\partial w} &= \frac{1}{3} (300 \cdot (150w + b - 200) + 400 \cdot (200w + b - 600) + 520 \cdot (260w + b - 500)) \\ &= \frac{1}{3} (260200w + 1220b - 560000) \end{aligned}$$

Similarly, we find  $\frac{\partial J}{\partial b}$ :

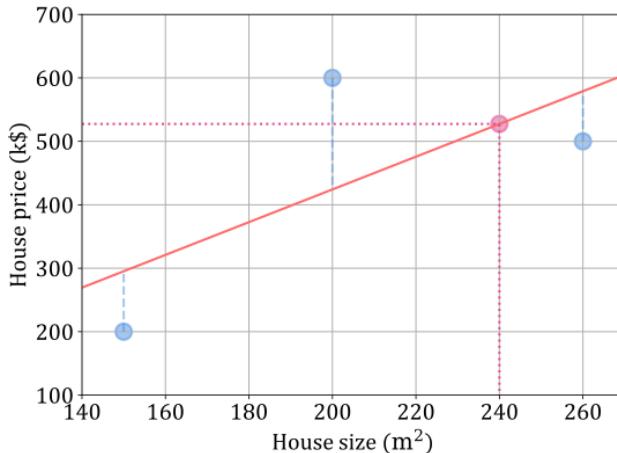
$$\begin{aligned} \frac{\partial J}{\partial b} &= \frac{1}{3} (2 \cdot (150w + b - 200) + 2 \cdot (200w + b - 600) + 2 \cdot (260w + b - 500)) \\ &= \frac{1}{3} (1220w + 6b - 2600) \end{aligned}$$

Setting the partial derivatives to 0, which is required to locate the optimal values, results in the following system of equations:

$$\begin{cases} \frac{1}{3} (260200w + 1220b - 560000) = 0 \\ \frac{1}{3} (1220w + 6b - 2600) = 0 \end{cases}$$

Simplifying the system and using substitution to solve for the variables gives the optimal values:  $w^* = 2.58$  and  $b^* = -91.76$ .

The resulting model  $f(x) = 2.58x - 91.76$  is shown in the plot below. It includes the three examples (blue dots), the model itself (red solid line), and a prediction for a new house with an area of  $240 \text{ m}^2$  (dotted orange lines).



A vertical blue dashed line shows the square root of the model's prediction error compared to the actual price.<sup>1</sup> Smaller errors mean the model **fits** the data better. The loss, which aggregates these errors, measures how well the model aligns with the dataset.

When we calculate the loss using the same dataset that trained the model, the result is called the **training loss**. The dataset used for training is referred to as the training dataset or **training set**. For our model, the training loss is defined by Equation 1.3. Now, we can use the learned parameter values to compute the loss for the training set:

$$\begin{aligned}
 J(2.58, -91.76) &= \frac{(2.58 \cdot 150 - 91.76 - 200)^2}{3} + \frac{(2.58 \cdot 200 - 91.76 - 600)^2}{3} \\
 &\quad + \frac{(2.58 \cdot 260 - 91.76 - 500)^2}{3} \\
 &= 15403.19.
 \end{aligned}$$

The square root of this value is approximately 124.1, indicating an average prediction error of around \$124,100. The interpretation of whether a loss value is high or low depends on the specific business context and comparative benchmarks. Neural networks and other non-linear models, which we explore later in this chapter, typically achieve lower loss values.

### 1.3. Four-Step Machine Learning Process

At this stage, you should clearly understand the four steps involved in supervised learning:

1. **Collect a dataset:** For example,  $(x_1, y_1) = (150, 200)$ ,  $(x_2, y_2) = (200, 600)$ , and  $(x_3, y_3) = (260, 500)$ .

---

<sup>1</sup> It's the square root of the error because our error, as defined in Equation 1.2, is the square of the difference between the predicted price and the real price of the house. It's common practice to take the square root of the mean squared error because it expresses the error in the same units as the target variable (price in this case). This makes it easier to interpret the error value.

2. **Define the model's structure:** For example,  $y = wx + b$ .
3. **Define the loss function:** Such as Equation 1.3.
4. **Minimize the loss:** Minimize the loss function on the dataset.

In our example, we minimized the loss manually by solving a system of two equations with two variables. This approach works for small systems. However, as models grow in complexity—such as large language models with billions of parameters—the manual approach becomes infeasible. Let's now introduce new concepts that will help us address this challenge.

## 1.4. Vector

To predict a house price, knowing its area alone isn't enough. Factors like the year of construction or the number of bedrooms and bathrooms also matter. Suppose we use two attributes: (1) area and (2) number of bedrooms. In this case, the input  $\mathbf{x}$  becomes a **feature vector**. This vector includes two **features**, also called **dimensions** or **components**:

$$\mathbf{x} \stackrel{\text{def}}{=} \begin{bmatrix} x^{(1)} \\ x^{(2)} \end{bmatrix}$$

In this book, vectors are represented with lowercase bold letters, such as  $\mathbf{x}$  or  $\mathbf{w}$ . For a given house  $\mathbf{x}$ ,  $x^{(1)}$  represents its size in square meters, and  $x^{(2)}$  represents the number of bedrooms. The **dimensionality** of the vector, or its **size**, refers to the number of components it contains. Here,  $\mathbf{x}$  has two components, so its dimensionality is 2.

A vector is usually represented as a column of numbers, called a **column vector**. However, in text, it is often written as its **transpose**,  $\mathbf{x}^\top$ . Transposing a column vector converts it into a **row vector**. For example,  $\mathbf{x}^\top \stackrel{\text{def}}{=} [x^{(1)}, x^{(2)}]$ .

With two features, our linear model needs three parameters: the weights  $w^{(1)}$  and  $w^{(2)}$ , and the bias  $b$ . The weights can be grouped into a vector:

$$\mathbf{w} \stackrel{\text{def}}{=} \begin{bmatrix} w^{(1)} \\ w^{(2)} \end{bmatrix}$$

The linear model can then be written compactly as:

$$y = \mathbf{w} \cdot \mathbf{x} + b, \tag{1.4}$$

where  $\mathbf{w} \cdot \mathbf{x}$  is a **dot product** of two vectors (also known as **scalar product**). It is defined as:

$$\mathbf{w} \cdot \mathbf{x} \stackrel{\text{def}}{=} \sum_{j=1}^D w^{(j)} x^{(j)}$$

The dot product combines two vectors of the same dimensionality to produce a **scalar**, a single number like 22, 0.67, or -10.5. Scalars in this book are denoted by italic lowercase or uppercase

letters, such as  $x$  or  $D$ . The expression  $\mathbf{w} \cdot \mathbf{x} + b$  generalizes the idea of a **linear transformation** to vectors.

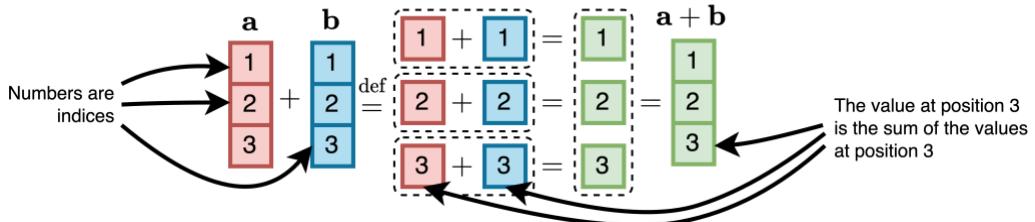
The equation uses **capital-sigma notation**, where  $D$  represents the dimensionality of the input, and  $j$  runs from 1 to  $D$ . For example, in the 2-dimensional house scenario,  $\sum_{j=1}^2 w^{(j)} x^{(j)} \stackrel{\text{def}}{=} w^{(1)}x^{(1)} + w^{(2)}x^{(2)}$ .

Although the capital-sigma notation suggests the dot product might be implemented as a loop, modern computers handle it much more efficiently. Optimized **linear algebra libraries** like **BLAS** and **cuBLAS** compute the dot product using low-level, highly optimized methods. These libraries leverage hardware acceleration and parallel processing, achieving speeds far beyond a simple manual loop.

The **sum of two vectors**  $\mathbf{a}$  and  $\mathbf{b}$ , both with the same dimensionality  $D$ , is defined as:

$$\mathbf{a} + \mathbf{b} \stackrel{\text{def}}{=} (a^{(1)} + b^{(1)}, a^{(2)} + b^{(2)}, \dots, a^{(D)} + b^{(D)})^\top$$

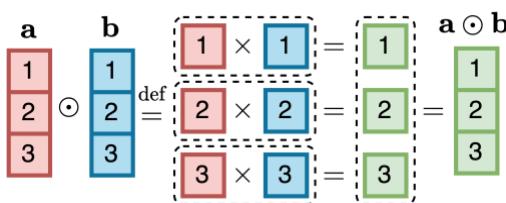
The calculation for a sum of two 3-dimensional vectors is illustrated below:<sup>2</sup>



The **element-wise product** of two vectors  $\mathbf{a}$  and  $\mathbf{b}$  of dimensionality  $D$ , is defined as:

$$\mathbf{a} \odot \mathbf{b} \stackrel{\text{def}}{=} (a^{(1)} \cdot b^{(1)}, a^{(2)} \cdot b^{(2)}, \dots, a^{(D)} \cdot b^{(D)})^\top$$

The computation of the element-wise product for two 3-dimensional vectors is shown below:



<sup>2</sup> In this chapter's illustrations, the numbers in the cells indicate the position of an element within an input or output matrix, or a vector. They do not represent actual values.

The **norm** of a vector  $\mathbf{x}$ , denoted  $\|\mathbf{x}\|$ , represents its **length** or **magnitude**. It is defined as the square root of the sum of the squares of its components:

$$\|\mathbf{x}\| \stackrel{\text{def}}{=} \sqrt{\sum_{j=1}^D (x^{(j)})^2}$$

For a 2-dimensional vector  $\mathbf{x}$ , the norm is:

$$\|\mathbf{x}\| = \sqrt{(x^{(1)})^2 + (x^{(2)})^2}$$

The cosine of the angle  $\theta$  between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  is defined as:

$$\cos(\theta) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \quad (1.5)$$

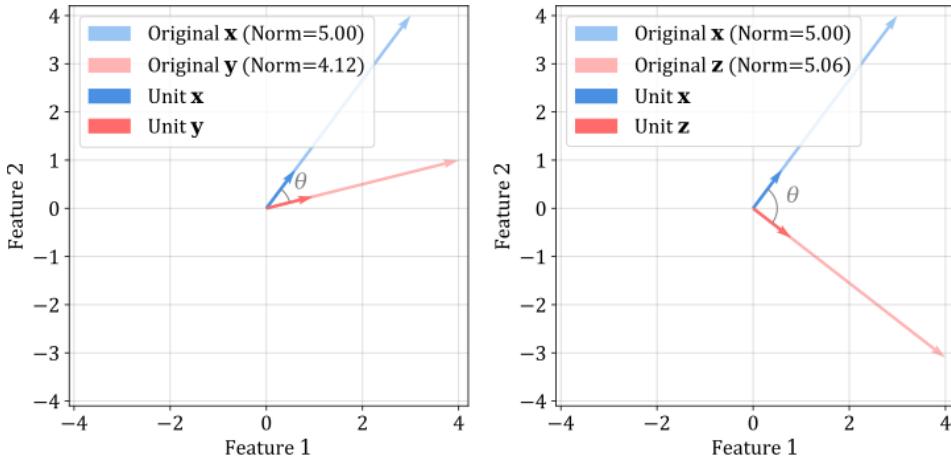
The cosine of the angle between two vectors quantifies their similarity. For instance, two houses with similar areas and bedroom counts will have a cosine similarity close to 1. **Cosine similarity** is widely used to compare words or documents represented as **embedding vectors**. This will be discussed further in Section 2.2.

A **zero vector** has all components equal to zero. A **unit vector** has a length of 1. To convert any non-zero vector  $\mathbf{x}$  into a unit vector  $\hat{\mathbf{x}}$ , you divide the vector by its norm:

$$\hat{\mathbf{x}} = \frac{\mathbf{x}}{\|\mathbf{x}\|}$$

Dividing a vector by a number results in a new vector where each component of the original vector is divided by that number.

A unit vector preserves the direction of the original vector but has a length of 1. The figure below demonstrates this with 2-dimensional examples. On the left, aligned vectors have  $\cos(\theta) = 0.78$ . On the right, nearly orthogonal vectors have  $\cos(\theta) = -0.02$ .



Unit vectors are valuable because their dot product equals the cosine of the angle between them, and computing dot products is efficient. When documents are represented as unit vectors, finding similar ones becomes fast by calculating the dot product between the query vector and document vectors. This is how vector search engines and libraries like Faiss, Qdrant, and Weaviate operate.

As dimensions increase, the number of parameters in a linear model becomes too large to solve manually. These models also face inherent limitations—they can only fit data that follows a straight line or its higher-dimensional analogues like planes and hyperplanes. (This problem is illustrated in the next section.)

In high-dimensional spaces, we cannot visually verify if data follows a linear pattern. Even if we could visualize beyond three dimensions, we would still need more flexible models to handle data that linear models cannot fit.

The next section explores non-linear models, with a focus on neural networks—the foundation for understanding large language models, which are a specialized neural network architecture.

## 1.5. Neural Network

A **neural network** differs from a linear model in two fundamental ways: (1) it applies fixed non-linear functions to the outputs of trainable linear functions, and (2) its structure is deeper, combining multiple functions hierarchically through layers. Let's illustrate these differences.

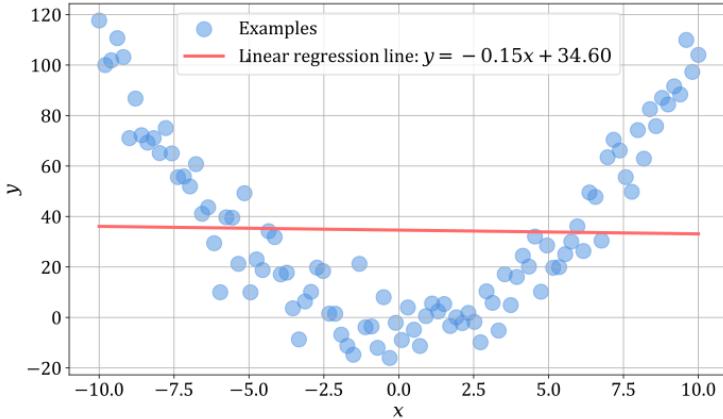
Linear models like  $wx + b$  or  $\mathbf{w} \cdot \mathbf{x} + b$  cannot solve many machine learning problems effectively. Even if we combine them into a **composite function**  $f_2(f_1(x))$ , a composite function of linear functions remains linear. This is straightforward to verify.

Let's define  $y_1 = f_1(x) \stackrel{\text{def}}{=} a_1x$  and  $y_2 = f_2(y_1) \stackrel{\text{def}}{=} a_2y_1$ . Here,  $f_2$  depends on  $f_1$ , making it a composite function. We can rewrite  $f_2$  as:

$$y_2 = a_2 y_1 = a_2(a_1 x) = (a_2 a_1)x$$

Since  $a_1$  and  $a_2$  are constants, we can define  $a_3 \stackrel{\text{def}}{=} a_1 a_2$ , so  $y_2 = a_3 x$ , which is linear.

A straight line often fails to capture patterns in one-dimensional data, as demonstrated when **linear regression** is applied to non-linear data:



To overcome this, we introduce non-linearity. For a 1D input, the model becomes:

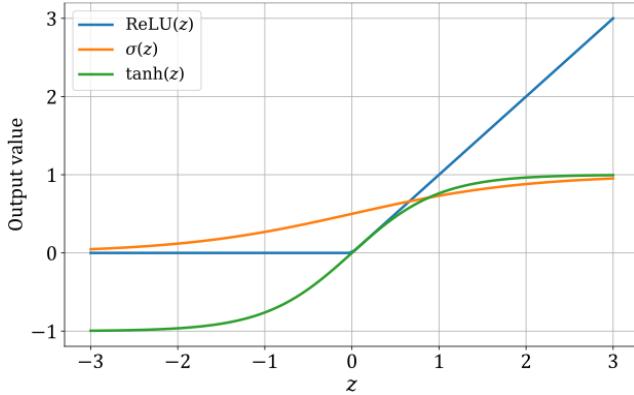
$$y = \phi(wx + b)$$

The function  $\phi$  is a fixed non-linear function, known as an **activation**. Common choices are:

- 1) **ReLU (rectified linear unit)**:  $\text{ReLU}(z) \stackrel{\text{def}}{=} \max(0, z)$ , which outputs non-negative values and is widely used in neural networks;
- 2) **Sigmoid**:  $\sigma(z) \stackrel{\text{def}}{=} \frac{1}{1+e^{-z}}$ , which outputs values between 0 and 1, making it suitable for **binary classification** (e.g., classifying spam emails as 1 and non-spam as 0);
- 3) **Tanh (hyperbolic tangent)**:  $\tanh(z) \stackrel{\text{def}}{=} \frac{e^z - e^{-z}}{e^z + e^{-z}}$ ; outputs values between  $-1$  and  $1$ .

In these equations,  $e$  denotes **Euler's number**, approximately 2.72.

These functions are widely used due to their mathematical properties, simplicity, and effectiveness in diverse applications. This is what they look like:

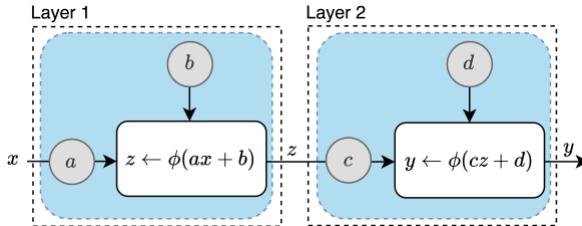


The structure  $\phi(wx + b)$  enables learning non-linear models but can't capture all non-linear curves. By nesting these functions, we build more expressive models. For instance, let  $f_1(x) \stackrel{\text{def}}{=} \phi(ax + b)$  and  $f_2(z) \stackrel{\text{def}}{=} \phi(cz + d)$ . A **composite model** combining  $f_1$  and  $f_2$  is:

$$y = f_2(f_1(x)) = \phi(c\phi(ax + b) + d)$$

Here, the input  $x$  is first transformed linearly using parameters  $a$  and  $b$ , then passed through the non-linear function  $\phi$ . The result is further transformed linearly with parameters  $c$  and  $d$ , followed by another application of  $\phi$ .

Below is the graph representation of the composite model  $y = f_2(f_1(x))$ :



A **computational graph** represents the structure of a model. The computational graph above shows two non-linear **units** (blue rectangles), often referred to as **artificial neurons**. Each unit contains two trainable parameters—a weight and a bias—represented by grey circles. The left arrow  $\leftarrow$  denotes that the value on the right is assigned to the variable on the left. This graph illustrates a basic neural network with two **layers**, each containing one unit. Most neural networks in practice are built with more layers and multiple units per layer.

Suppose we have a two-dimensional input. The **input layer** contains three units, while the **output layer** has a single unit. The network's structure appears as follows:

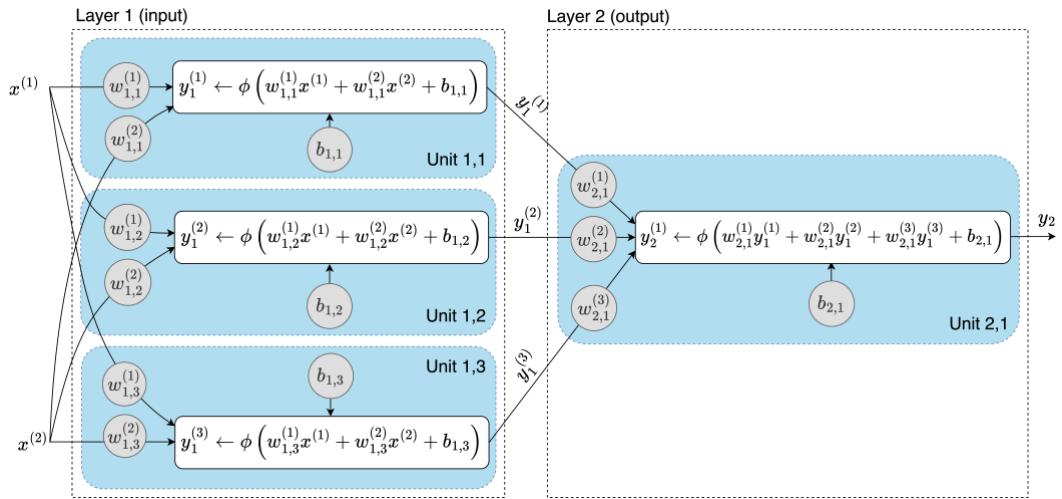


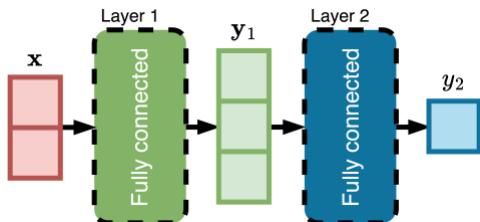
Figure 1.1: A neural network with two layers.

This structure represents a **feedforward neural network (FNN)**, where information flows in one direction—left to right—without loops. When units in each layer connect to all units in the subsequent layer, as shown above, we call it a **multilayer perceptron (MLP)**. A layer where each unit connects to all units in both adjacent layers is termed a **fully connected layer**, or **dense layer**.

In Chapter 3, we will explore recurrent neural networks (RNNs). Unlike FNNs, RNNs have loops, where outputs from a layer are used as inputs to the same layer.

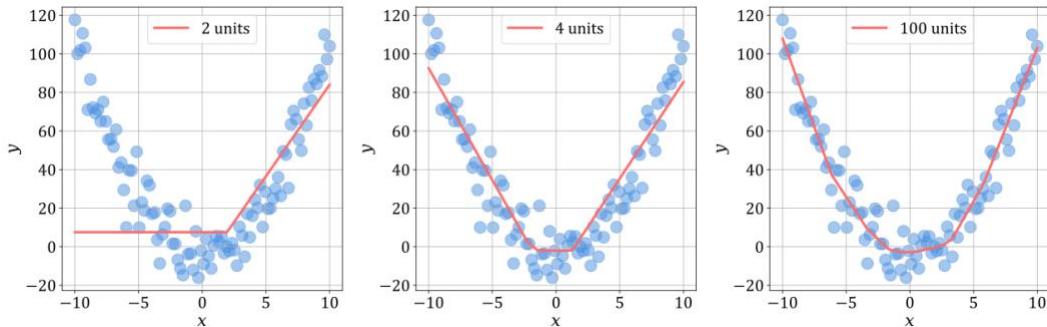
**Convolutional neural networks (CNN)** are feedforward neural networks with convolutional layers that are not fully connected. While initially designed for image processing, they are effective for tasks like document classification in text data. To learn more about CNNs refer to the additional materials in the book's wiki.

To simplify diagrams, individual neural units can be replaced with squares. Using this approach, the above network can be represented more compactly as follows:



If you think this simple model is too weak, look at the figure below. It contains three plots demonstrating how increasing model size improves performance. The left plot shows a model with

2 units: one input, one output, and ReLU activations. The middle plot is a model with 4 units: three inputs and one output. The right plot shows a much larger model with 100 units:



The ReLU activation function, despite its simplicity, was a breakthrough in machine learning. Neural networks before 2012 relied on smooth activations like tanh and sigmoid, which made training deep models increasingly difficult. We will return to this subject in Chapter 4 on the Transformer neural network architecture.

Increasing the number of parameters helps the model approximate the data more accurately. Experiments consistently show that adding more units per layer or increasing the number of layers in a neural network improves its capacity to fit high-dimensional datasets, such as natural language, voice, sound, image, and video data.

## 1.6. Matrix

Neural networks can handle high-dimensional datasets but require substantial memory and computation. Calculating a layer's transformation naïvely would involve iterating over thousands of parameters per unit across thousands of units and dozens of layers, which is both slow and resource-intensive. Using **matrices** makes the computations more efficient.

A **matrix** is a two-dimensional array of numbers arranged into rows and columns, which generalizes the concept of vectors to higher dimensionalities. Formally, a matrix  $\mathbf{A}$  with  $m$  rows and  $n$  columns is written as:

$$\mathbf{A} \stackrel{\text{def}}{=} \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix}$$

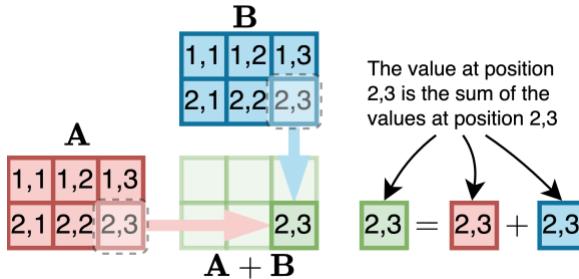
Here,  $a_{i,j}$  represents the element in the  $i$ -th row and  $j$ -th column of the matrix. The dimensions of the matrix are expressed as  $m \times n$  (read as "m by n").

Matrices are fundamental in machine learning. They compactly represent data and weights and enable efficient computation through operations such as addition, multiplication, and transposition. In this book, matrices are represented with uppercase bold letters, such as  $\mathbf{X}$  or  $\mathbf{W}$ .

The **sum of two matrices A and B** of the same dimensionality is defined element-wise as:

$$(\mathbf{A} + \mathbf{B})_{i,j} \stackrel{\text{def}}{=} a_{i,j} + b_{i,j}$$

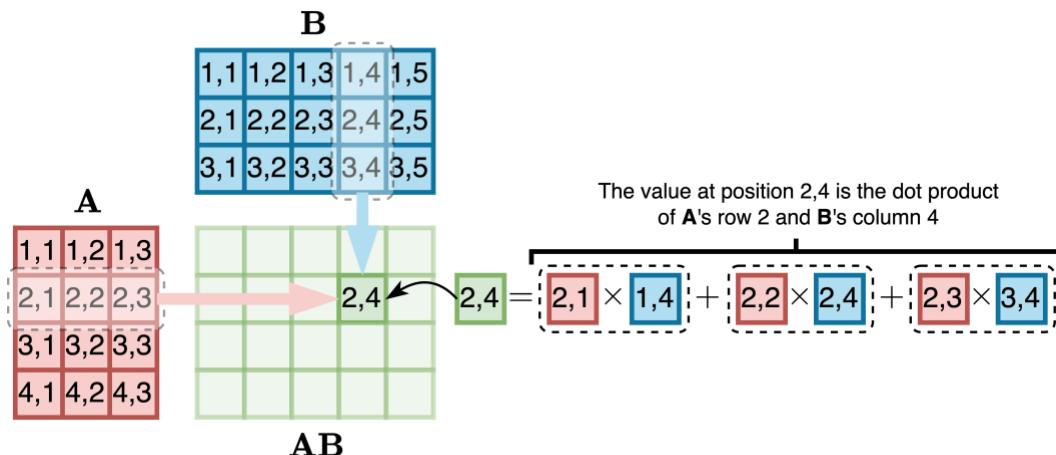
For example, for two  $2 \times 3$  matrices **A** and **B**, the addition works like this:



The **product of a matrices A** with dimensions  $m \times n$  and **B** with dimensions  $n \times p$  is a matrix **C** with dimensions  $m \times p$  such that the value in row  $i$  and column  $k$  is given by:

$$(\mathbf{C})_{i,k} = \sum_{j=1}^n a_{i,j} b_{j,k}$$

For example, for a  $4 \times 3$  matrix **A** and a  $3 \times 5$  matrix **B**, the product is a  $4 \times 5$  matrix:



**Transposing a matrix A** swaps its rows and columns, resulting in  $\mathbf{A}^\top$ , where:

$$(\mathbf{A}^\top)_{i,j} = a_{j,i}$$

For example, for a  $2 \times 3$  matrix  $\mathbf{A}$ , its transpose  $\mathbf{A}^T$  look like this:

$$\begin{array}{c} \mathbf{A} \\ \begin{array}{|c|c|c|} \hline 1 & 3 & 5 \\ \hline 2 & 4 & 6 \\ \hline \end{array} \end{array} \quad \begin{array}{c} \mathbf{A}^T \\ \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline 5 & 6 \\ \hline \end{array} \end{array}$$

**Matrix-vector multiplication** is a special case of matrix multiplication. When an  $m \times n$  matrix  $\mathbf{A}$  is multiplied by a vector  $\mathbf{x}$  of size  $n$ , the result is a vector  $\mathbf{y} = \mathbf{Ax}$  with  $m$  components.

Each element  $y_i$  of the resulting vector  $\mathbf{y}$  is computed as:

$$y_i = \sum_{j=1}^n a_{i,j} x^{(j)}$$

For example, a  $4 \times 3$  matrix  $\mathbf{A}$  multiplied by a 3D vector  $\mathbf{x}$  produces a 4-dimensional vector:

$$\begin{array}{c} \mathbf{A} \\ \begin{array}{|c|c|c|} \hline 1,1 & 1,2 & 1,3 \\ \hline 2,1 & 2,2 & 2,3 \\ \hline 3,1 & 3,2 & 3,3 \\ \hline 4,1 & 4,2 & 4,3 \\ \hline \end{array} \end{array} \times \begin{array}{c} \mathbf{x} \\ \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 3 \\ \hline \end{array} \end{array} = \begin{array}{l} \boxed{1,1} \times \boxed{1} + \boxed{1,2} \times \boxed{2} + \boxed{1,3} \times \boxed{3} = \boxed{1} \\ \boxed{2,1} \times \boxed{1} + \boxed{2,2} \times \boxed{2} + \boxed{2,3} \times \boxed{3} = \boxed{2} \\ \boxed{3,1} \times \boxed{1} + \boxed{3,2} \times \boxed{2} + \boxed{3,3} \times \boxed{3} = \boxed{3} \\ \boxed{4,1} \times \boxed{1} + \boxed{4,2} \times \boxed{2} + \boxed{4,3} \times \boxed{3} = \boxed{4} \end{array} = \begin{array}{c} \mathbf{Ax} \\ \begin{array}{|c|c|c|c|} \hline 1 & 2 & 3 & 4 \\ \hline \end{array} \end{array}$$

The weights and biases in fully connected layers of neural networks can be compactly represented using matrices and vectors, enabling the use of highly optimized linear algebra libraries. As a result, matrix operations form the backbone of neural network training and inference.

Let's express the model in Figure 1.1 using matrix notation. Let  $\mathbf{x}$  be the 2D input feature vector. For the first layer, the weights and biases are represented as a  $3 \times 2$  matrix  $\mathbf{W}_1$  and a 3D vector  $\mathbf{b}_1$ , respectively. The 3D output  $\mathbf{y}_1$  of the first layer is given by:

$$\mathbf{y}_1 = \phi(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \quad (1.6)$$

The second layer also uses a weight matrix and a bias. The output  $\mathbf{y}_2$  of the second layer is computed using the output  $\mathbf{y}_1$  from the first layer. The weight matrix for the second layer is a  $1 \times 3$  matrix  $\mathbf{W}_2$ . The bias for the second layer is a scalar  $b_{2,1}$ . The model output corresponds to the output of the second layer:

$$\mathbf{y}_2 = \phi(\mathbf{W}_2 \mathbf{y}_1 + b_{2,1}) \quad (1.7)$$

Equation 1.6 and Equation 1.7 capture the operations from input to output in the neural network, with each layer's output serving as the input for the next.

## 1.7. Gradient Descent

Neural networks are typically large and composed of non-linear functions, which makes solving for the minimum of the loss function analytically infeasible. Instead, the gradient descent algorithm is widely used to minimize the loss, including in large language models.

Consider a practical example: **binary classification**. This task assigns input data to one of two classes, like deciding if an email is spam or not, or detecting whether a website connection request is a DDoS attack.

Our training dataset  $\mathcal{D}$  is  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , where  $\mathbf{x}_i$  are vectors of input features, and  $y_i$  are the labels. Each  $y_i$ , indexed from 1 to  $N$ , takes a value of 0 for "not spam" or 1 for "spam." A well-trained model should output  $\tilde{y}$  close to 1 for spam inputs  $\mathbf{x}$  and close to 0 for non-spam inputs. We can define the model as follows:

$$y = \sigma(\mathbf{w} \cdot \mathbf{x} + b), \quad (1.8)$$

where  $\mathbf{x} = [x^{(j)}]_{j=1}^D$  and  $\mathbf{w} = [w^{(j)}]_{j=1}^D$  are  $D$ -dimensional vectors,  $b$  is a scalar, and  $\sigma$  is the **sigmoid** defined in Section 1.5.

This model, called **logistic regression**, is commonly used for binary classification tasks. Unlike **linear regression**, which produces outputs ranging from  $-\infty$  to  $\infty$ , logistic regression always outputs values between 0 and 1. It can serve either as a standalone model or as the output layer in a larger neural network.

Despite being over 80 years old, logistic regression remains one of the most widely used algorithms in production machine learning systems.

A common choice for the **loss function** in this case is **binary cross-entropy**, also called **logistic loss**. For a single example  $i$ , the binary cross-entropy loss is defined as:

$$\text{loss}(\tilde{y}_i, y_i) \stackrel{\text{def}}{=} -[y_i \log(\tilde{y}_i) + (1 - y_i) \log(1 - \tilde{y}_i)] \quad (1.9)$$

In this equation,  $y_i$  represents the actual label of the  $i$ -th example in the dataset, and  $\tilde{y}_i$  is the **prediction score**, a value between 0 and 1 that the model outputs for input vector  $\mathbf{x}_i$ . The function  $\log$  denotes the **natural logarithm**.

Loss functions are usually designed to penalize incorrect predictions while rewarding accurate ones. To see why logistic loss works for logistic regression, consider two extreme cases:

1. **Perfect prediction**, when  $y_i = 0$  and  $\tilde{y}_i = 0$ :

$$\text{loss}(0,0) = -[0 \cdot \log(0) + (1 - 0) \cdot \log(1 - 0)] = -\log(1) = 0$$

Here, the loss is zero which is good because the prediction matches the label perfectly.

## 2. Opposite prediction, when $y_i = 0$ and $\tilde{y}_i = 1$ :

$$\text{loss}(1,0) = -[0 \cdot \log(1) + (1-0) \cdot \log(1-1)] = -\log(0)$$

The logarithm of 0 is undefined, and as  $a$  approaches 0,  $-\log(a)$  approaches infinity, representing a severe loss for completely wrong predictions. However, since  $\tilde{y}_i$ , the output of the sigmoid function, always remains strictly between 0 and 1, the loss stays finite.

For an entire dataset  $\mathcal{D}$ , the loss is given by the average loss for all examples in the dataset:

$$\text{loss}_{\mathcal{D}} \stackrel{\text{def}}{=} -\frac{1}{N} \sum_{i=1}^N [y_i \log(\tilde{y}_i) + (1-y_i) \log(1-\tilde{y}_i)] \quad (1.10)$$

To simplify the gradient descent derivation, we'll stick to a single example,  $i$ , and rewrite the equation by substituting the prediction score  $\tilde{y}_i$  with the model's expression for it:

$$\text{loss}(\tilde{y}_i, y_i) = -[y_i \log(\sigma(z_i)) + (1-y_i) \log(1-\sigma(z_i))], \text{ where } z_i = \mathbf{w} \cdot \mathbf{x}_i + b$$

To minimize  $\text{loss}(\tilde{y}_i, y_i)$ , we calculate the partial derivatives with respect to each weight  $w^{(j)}$  and the bias  $b$ . We will use the **chain rule** because we have a **composition** of three functions:

- **Function 1:**  $z_i \stackrel{\text{def}}{=} \mathbf{w} \cdot \mathbf{x}_i + b$ , a linear function involving the weights  $\mathbf{w}$  and the bias  $b$ ;
- **Function 2:**  $\tilde{y}_i = \sigma(z_i) \stackrel{\text{def}}{=} \frac{1}{1+e^{-z_i}}$ , the sigmoid function applied to  $z_i$ ;
- **Function 3:**  $\text{loss}(\tilde{y}_i, y_i)$ , as defined in Equation 1.9, which depends on  $\tilde{y}_i$ .

Notice that  $\mathbf{x}_i$  and  $y_i$  are given:  $\mathbf{x}_i$  is the feature vector for example  $i$ , and  $y_i \in \{0,1\}$  is its label. The notation  $y_i \in \{0,1\}$  means that  $y_i$  belongs to the set  $\{0,1\}$  and, in this case, indicates that  $y_i$  can only be 0 or 1.

Let's denote  $\text{loss}(\tilde{y}_i, y_i)$  as  $l_i$ . For weights  $w^{(j)}$ , the application of the chain rule gives us:

$$\frac{\partial l_i}{\partial w^{(j)}} = \frac{\partial l_i}{\partial \tilde{y}_i} \cdot \frac{\partial \tilde{y}_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w^{(j)}} = (\tilde{y}_i - y_i) \cdot x_i^{(j)}$$

For the bias  $b$ , we have:

$$\frac{\partial l_i}{\partial b} = \frac{\partial l_i}{\partial \tilde{y}_i} \cdot \frac{\partial \tilde{y}_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial b} = \tilde{y}_i - y_i$$

This is where the beauty of machine learning math truly shines: the activation function—sigmoid—and loss function—cross-entropy—both arise from  $e$ , Euler's number. Their functional properties serve distinct purposes: sigmoid ranges between 0 and 1, ideal for binary classification, while cross-entropy spans from 0 to  $\infty$ , perfect

as a penalty. When combined, the exponential and logarithmic components elegantly cancel, yielding a linear function—prized for its computational simplicity and numerical stability. The book's wiki provides the full derivation.

The partial derivatives with respect to  $w^{(j)}$  and  $b$  for a single example  $(\mathbf{x}_i, y_i)$  can be extended to the entire dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  by summing the contributions from all examples and averaging them. This follows from the **sum rule** and the **constant multiple rule** of differentiation:

$$\begin{aligned}\frac{\partial \text{loss}}{\partial w^{(j)}} &= \frac{1}{N} \sum_{i=1}^N [(\tilde{y}_i - y_i) \cdot x_i^{(j)}] \\ \frac{\partial \text{loss}}{\partial b} &= \frac{1}{N} \sum_{i=1}^N [\tilde{y}_i - y_i]\end{aligned}\tag{1.11}$$

Averaging the losses for individual examples ensures that each example contributes equally to the overall loss, regardless of the total number of examples.

The **gradient** is a vector that contains all the partial derivatives. The gradient of the loss function, denoted as  $\nabla \text{loss}$ , is defined as follows:

$$\nabla \text{loss} \stackrel{\text{def}}{=} \left( \frac{\partial \text{loss}}{\partial w^{(1)}}, \frac{\partial \text{loss}}{\partial w^{(2)}}, \dots, \frac{\partial \text{loss}}{\partial w^{(D)}}, \frac{\partial \text{loss}}{\partial b} \right)$$

If a gradient's component is positive, this means that increasing the corresponding parameter will increase the loss. Therefore, to minimize the loss, we should decrease that parameter.

The **gradient descent** algorithm uses the gradient of the loss function to iteratively update the weights and bias, aiming to minimize the loss function. Here's how it operates:

0. **Initialize parameters:** Start with random values of parameters  $w^{(j)}$  and  $b$ .
1. **Compute the predictions:** For each training example  $(\mathbf{x}_i, y_i)$ , compute the predicted value  $\tilde{y}_i$  using the model:

$$\tilde{y}_i \leftarrow \sigma(\mathbf{w} \cdot \mathbf{x}_i + b)$$

2. **Compute the gradient:** Calculate the partial derivatives of the loss function with respect to each weight  $w^{(j)}$  and the bias  $b$  using Equation 1.11.
3. **Update the weights and bias:** Adjust the weights and bias in the direction that decreases the loss function. This adjustment involves taking a small step in the opposite direction of the gradient. The step size is controlled by the learning rate  $\eta$  (explained below):

$$w^{(j)} \leftarrow w^{(j)} - \eta \frac{\partial \text{loss}}{\partial w^{(j)}}$$

$$b \leftarrow b - \eta \frac{\partial \text{loss}}{\partial b}$$

4. **Calculate the loss:** Calculate the logistic loss by substituting the updated values of  $w^{(j)}$  and  $b$  into Equation 1.10.
5. **Continue the iterative process:** Repeat steps 1-4 for a set number of **iterations** (also called **steps**) or until the loss value converges to a minimum.

Here's a bit more detail to clarify the steps:

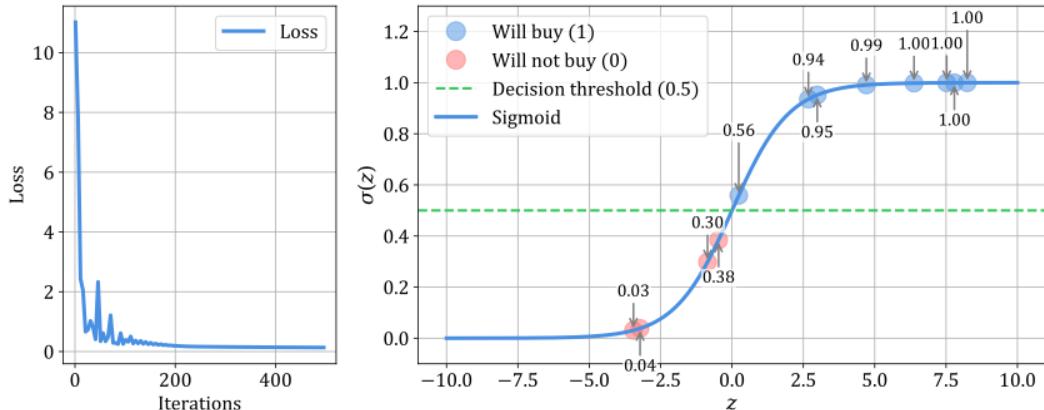
- Gradients are subtracted from parameters because they point in the direction of steepest ascent in the loss function. Since our goal is to minimize loss, we move in the opposite direction—hence, the subtraction.
- The **learning rate**  $\eta$  is a positive value close to 0 and serves as a **hyperparameter**,—not learned by the model but set manually. It controls the step size of each update, and finding its optimal value requires experimentation.
- **Convergence** occurs when subsequent iterations yield minimal decreases in loss. The learning rate  $\eta$  is crucial here: too small, and progress crawls; too large, and we risk overshooting the minimum or even seeing the loss increase rather than decrease. Choosing an appropriate  $\eta$  is therefore essential for effective gradient descent.

Let's illustrate the process with a simple dataset of 12 examples:

$$\left\{ ((22, 25), 0), ((25, 35), 0), ((47, 80), 1), ((52, 95), 1), ((46, 82), 1), ((56, 90), 1), ((23, 27), 0), ((30, 50), 1), ((40, 60), 1), ((39, 57), 0), ((53, 95), 1), ((48, 88), 1) \right\}$$

In this dataset,  $x_i$  contains two features: age (in years) and income (in thousands of dollars). The objective is to predict whether a person will buy a product, with label  $y_i$  being either 0 (will not buy) or 1 (will buy).

The loss evolution across gradient descent steps and the resulting trained model are shown in the figure below:



The left plot shows the loss decreasing steadily during gradient descent optimization. The right plot displays the trained model's sigmoid function, with training examples positioned by their z-values ( $z_i = \mathbf{w}^* \cdot \mathbf{x}_i + b^*$ ), where  $\mathbf{w}^*$  and  $b^*$  are the learned weights and bias.

The 0.5 threshold was chosen based on the plot's clear separation: all "will-buy" examples (blue dots) lie above it, while all "will-not-buy" examples (red dots) fall below. For new inputs  $\mathbf{x}$ , predict using  $\tilde{y} = \sigma(\mathbf{w}^* \cdot \mathbf{x} + b^*)$ . If  $\tilde{y} < 0.5$ , predict "will not buy;" otherwise, "will buy."

## 1.8. Automatic Differentiation

Gradient descent optimizes model parameters but requires partial derivative equations. Until now, these derivatives had to be calculated by hand for each model. As models grow more complex, particularly in neural networks with multiple layers, manual derivation becomes impractical.

This is where **automatic differentiation** (or **autograd**) comes in. Built into machine learning frameworks like PyTorch and TensorFlow, this feature computes partial derivatives directly from model-defining Python code. This eliminates manual derivation, even for the most sophisticated models.

Modern automatic differentiation systems can handle derivatives for millions of variables efficiently. Manual computation of these derivatives would be unfeasible—writing the equations alone could take years.

To use gradient descent in PyTorch, first install it with `pip3` like this:

```
$ pip3 install torch
```

Now that PyTorch is installed, let's import the dependencies:

```
import torch
import torch.nn as nn
import torch.optim as optim
```

The `torch.nn` module contains building blocks for creating models. When you use these components, PyTorch automatically handles derivative calculations. For optimization algorithms like gradient descent, the `torch.optim` module has what you need. Here's how to implement logistic regression in PyTorch:

```
model = nn.Sequential(
    nn.Linear(n_inputs, n_outputs), ❶
    nn.Sigmoid() ❷
)
```

Our model leverages PyTorch's **sequential API**, which is well-suited for simple feedforward neural networks where data flows sequentially through layers. Each layer's output naturally becomes the input for the subsequent layer. The more versatile **module API**, which we'll cover in the next chapter, enables the creation of models with multiple inputs, outputs, or loops.

The input layer, defined in line ❶ using `nn.Linear`, has input dimensionality (`n_inputs`) matching the size of our feature vector `x`, while the output dimensionality (`n_outputs`) determines the layer's unit count. For our buy/no-buy **classifier**—a model assigning classes to inputs—we set `n_inputs` to 2 since  $\mathbf{x} = [x^{(1)}, x^{(2)}]^\top$ . With the output `z` being scalar, `n_outputs` becomes 1. Line ❷ transforms `z` through the sigmoid function to produce the output score.

We then proceed to define our dataset, create the model instance, establish the binary cross-entropy loss function, and set up the gradient descent algorithm:

```
inputs = torch.tensor([
    [22, 25], [25, 35], [47, 80], [52, 95], [46, 82], [56, 90],
    [23, 27], [30, 50], [40, 60], [39, 57], [53, 95], [48, 88]
], dtype=torch.float32) ❶

labels = torch.tensor([
    [0], [0], [1], [1], [1], [0], [1], [1], [0], [1], [1]
], dtype=torch.float32) ❷

model = nn.Sequential(
    nn.Linear(inputs.shape[1], 1),
    nn.Sigmoid()
)
optimizer = optim.SGD(model.parameters(), lr=0.001)
criterion = nn.BCELoss() # binary cross-entropy loss
```

In the above code block, we defined `inputs` and `labels`. The `inputs` form a matrix with 12 rows and 2 columns, while the `labels` are a vector with 12 components. The `shape` attribute of the `inputs` tensor return its dimensionality:

```
>>> inputs.shape
torch.Size([12, 2])
```

**Tensors** are PyTorch's core data structures—multi-dimensional arrays optimized for computation on both CPU and GPU. Supporting automatic differentiation and flexible data reshaping, tensors form the foundation for neural network operations. In our example, the `inputs` tensor contains 12 examples with 2 features each, while the `labels` tensor holds 12 examples with single labels. Following standard convention, examples are arranged in rows and their features in columns.

If you're not familiar with tensors, there's an introductory chapter on tensors available on the book's wiki.

When creating tensors in PyTorch, specifying `dtype=torch.float32` in lines ❶ and ❷ sets 32-bit floating-point precision explicitly. This precision setting is essential for neural network computations, including weight adjustments, activation functions, and gradient calculations.

The 32-bit floating-point precision is not the only option for neural networks. **Quantization**, an advanced technique that uses lower-precision data types like 16-bit or 8-bit floats and integers, helps reduce model size and improve computational efficiency. For more information, refer to resources on model optimization and deployment available on the book's wiki.

The `optim.SGD` class implements gradient descent by taking a list of model parameters and learning rate as inputs.<sup>3</sup> Since our model inherits from `nn.Module`, we can access all trainable parameters through its `parameters` method.

PyTorch provides the **binary cross-entropy** loss function through `nn.BCELoss()`.

Now, we have everything we need to start the training loop:

```
for step in range(500):
    optimizer.zero_grad() ❶
    loss = criterion(model(inputs), labels) ❷
    loss.backward() ❸
    optimizer.step() ❹
```

Line ❷ calculates the binary cross-entropy loss (Equation 1.10) by evaluating model predictions against training labels. Line ❸ then uses backpropagation to compute the gradient of this loss with respect to the model parameters.

**Backpropagation** applies differentiation rules, particularly the chain rule, to compute gradients through deep composite functions. This algorithm forms the backbone of neural network training. When PyTorch operates on tensors, it builds a computational graph as shown in Figure 1.1 from Section 1.5. This graph tracks all operations performed on the tensors. The `loss.backward()` call prompts PyTorch to traverse this graph and compute gradients via the chain rule, eliminating the need for manual gradient derivation and implementation.

The flow of data from input to output through the computational graph constitutes the **forward pass**, while the computation of gradients from output to input through backpropagation represents the **backward pass**.

PyTorch accumulates gradients in the `.grad` attribute of parameters like weights and biases. While this feature enables multiple gradient computations before parameter updates—useful for recurrent neural networks (covered in Section 3)—our implementation doesn't require gradient accumulation. Line ❶ therefore clears the gradients at each step's beginning.

---

<sup>3</sup> While 0.001 is a common default learning rate, optimal values vary by problem and dataset. Finding the best rate involves systematically testing different values and comparing model performance.

Finally, in line ④, parameter values are updated by subtracting the product of the learning rate and the loss function's partial derivatives, completing step 3 of the gradient descent algorithm discussed earlier.

One of automatic differentiation's key advantages is its flexibility with model switching—as long as you're using PyTorch's components, you can readily swap between different architectures. For instance, you could replace logistic regression with a basic two-layer FNN, defined through the sequential API:

```
model = nn.Sequential(  
    nn.Linear(features.shape[1], 100),  
    nn.Sigmoid(),  
    nn.Linear(100, labels.shape[1]),  
    nn.Sigmoid()  
)
```

In this setup, each of the 100 units in the first layer contains 2 weights and 1 bias, while the output layer's single unit has 100 weights and 1 bias. The automatic differentiation system handles gradient computation internally, so the remaining code stays unchanged.

The next chapter examines how to represent and process text data, beginning with fundamental techniques for converting documents into numerical representations like bag-of-words and word embeddings, followed by exploring the count-based language modeling approach.

## Chapter 2. Language Modeling Basics

To understand language modeling, we first need to solve how to represent language in a way computers can handle. This chapter explores the transformation of documents and words into numerical formats for computation. We'll then define language models and examine our first architecture: count-based language models. The chapter concludes with methods for evaluating language models.

Let's begin with one of the oldest yet effective techniques for converting text into usable data for machine learning: bag of words.

### 2.1. Bag of Words

Suppose you have a collection of documents and want to predict the main topic of each one. When topics are defined in advance, this task is called **classification**. With only two possible topics, it's known as **binary classification**, as explained in Section 1.7. With more than two topics, we have **multiclass classification**.

In multiclass classification, the dataset consists of pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , where  $y_i \in \{1, \dots, C\}$ ,  $N$  represents the number of examples, and  $C$  denotes the number of possible classes. Each  $\mathbf{x}_i$  could be a text document, with  $y_i$  being an integer indicating its topic—for example, 1 for “music,” 2 for “science,” or 3 for “cinema.”

Machines don't process text like humans. To train a machine learning classifier, we first convert documents into numbers. Each document becomes a **feature vector**, where each feature is a **scalar** value.

A common and effective approach to convert a collection of documents into feature vectors is the **bag of words (BoW)**. Here's how it works for a collection of 10 documents as an example:

ID	Text
1	Movies are fun for everyone.
2	Watching movies is great fun.
3	Enjoy a great movie today.
4	Research is interesting and important.
5	Learning math is very important.
6	Science discovery is interesting.
7	Rock is great to listen to.
8	Listen to music for fun.
9	Music is fun for everyone.
10	Listen to folk music!

A collection of text documents used in machine learning is called a **corpus**. The bag of words method applied to a corpus involves two key steps:

1. **Create a vocabulary:** List all unique words in the corpus to create the **vocabulary**.
2. **Vectorize documents:** Convert each document into a feature vector, where each dimension represents a word from the vocabulary. The value indicates the word's presence, absence, or frequency in the document.

For the 10-document corpus, the vocabulary is built by listing all unique words in alphabetical order. This involves removing punctuation, converting words to lowercase, and eliminating duplicates. After processing, we get:

```
vocabulary = ["a", "and", "are", "discovery", "enjoy", "everyone", "folk",
  "for", "fun", "great", "important", "interesting", "is", "learning", "liste
n", "math", "movie", "movies", "music", "research", "rock", "science", "to",
  "today", "very", "watching"]
```

Splitting a document into small indivisible parts is called **tokenization**, and each part is a **token**. There are different ways to tokenize. We tokenized our 10-document corpus by words. Sometimes, it's useful to break words into smaller units, called **subwords** to keep the vocabulary size manageable. For instance, instead of including “interesting” in the vocabulary, we might split it into “interest” and “-ing.” A common method for subword tokenization, which we'll cover later in this chapter, is byte-pair encoding. The choice of tokenization method depends on the language, dataset, and model.

A count of all English word **surface forms**—like *do*, *does*, *doing*, and *did*—reveals several million possibilities. Languages with more complex morphology have even greater numbers. A Finnish noun alone can take 2,000–3,000 different forms to express various case and number combinations. Using subwords offers a practical solution, as storing every surface form in the vocabulary would consume excessive memory and computational resources.

Words are a type of token, so “token” and “word” are often used interchangeably as the smallest indivisible units of a document. When a distinction is important, context will make it clear. While the bag-of-words approach can handle both words and subwords, it was originally designed for words—hence the name.

Feature vectors can be organized into a **document-term matrix**. Here, rows represent documents, and columns represent tokens. Below is a partial document-term matrix for a 10-document corpus. It includes only a subset of tokens to fit within the page width:

Doc	a	and	...	fun	...	listen	math	...	science	...	watching
1	0	0	...	1	...	0	0	...	0	...	0
2	0	0	...	1	...	0	0	...	0	...	1
3	1	0	...	0	...	0	0	...	0	...	0
4	0	1	...	0	...	0	0	...	0	...	0
5	0	0	...	0	...	0	1	...	0	...	0
6	0	0	...	0	...	0	0	...	1	...	0

Doc	a	and	...	fun	...	listen	math	...	science	...	watching
7	0	0	...	0	...	1	0	...	0	...	0
8	0	0	...	1	...	1	0	...	0	...	0
9	0	0	...	1	...	0	0	...	0	...	0
10	0	0	...	0	...	1	0	...	0	...	0

In the document-term matrix above, a value of 1 means the token appears in the document, while 0 means it does not. For instance, the feature vector  $\mathbf{x}_2$  for document 2 ("Watching movies is great fun.") is:

$$\mathbf{x}_2 = [0,0,0,0,0,0,0,1,1,0,0,1,0,0,0,0,1,0,0,0,0,0,0,0,1]^T.$$

In natural languages, word frequencies follow **Zipf's Law**. This law states that a word's frequency is inversely proportional to its rank in a frequency table—for instance, the second most frequent word appears half as often as the most frequent one. As a result, most words in a document-term matrix are rare, creating a **sparse** matrix with numerous zeros.

A neural network can be trained to predict a document's topic using these feature vectors. Let's do that. The first step is to assign labels to the documents, a process known as **labeling**. Labeling can be done manually or assisted by an algorithm. When algorithms are used, human validation is often needed to confirm accuracy. Here, we will manually label the documents by reading each one and choosing the most suitable topic from the three options.

Doc	Text	Class ID	Class Name
1	Movies are fun for everyone.	1	Cinema
2	Watching movies is great fun.	1	Cinema
3	Enjoy a great movie today.	1	Cinema
4	Research is interesting and important.	3	Science
5	Learning math is very important.	3	Science
6	Science discovery is interesting.	3	Science
7	Rock is great to listen to.	2	Music
8	Listen to music for fun.	2	Music
9	Music is fun for everyone.	2	Music
10	Listen to folk music!	2	Music

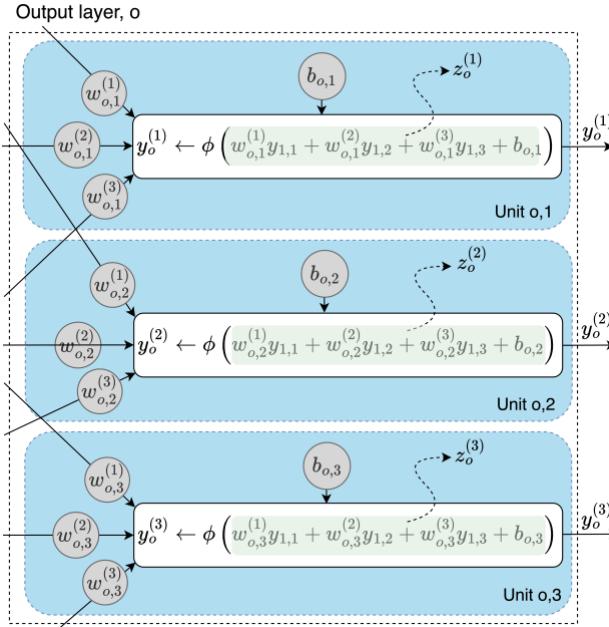
Advanced chat LMs enable highly accurate automated document labeling through a panel of expert models. Using three LLMs, when two or more assign the same label to a document, that label is adopted. If all three disagree, either a human can decide or a fourth model can break the tie. In many business contexts, manual labeling is becoming obsolete, as LLMs offer faster and more reliable labeling.

We have three classes: 1 for cinema, 2 for music, and 3 for science.<sup>4</sup> While binary classifiers typically use the **sigmoid** activation function with the **binary cross-entropy** loss, as discussed in Section 1.7, tasks involving three or more classes generally employ the softmax activation function paired with the cross-entropy loss.

The **softmax** function is defined as:

$$\text{softmax}(\mathbf{z}, k) \stackrel{\text{def}}{=} \frac{e^{z^{(k)}}}{\sum_{j=1}^D e^{z^{(j)}}}$$

Here,  $\mathbf{z}$  is a  $D$ -dimensional vector of logits,  $k$  is the index for which the softmax is computed, and  $e$  is **Euler's number**. **Logits** are the raw outputs of a neural network, prior to applying an activation function, as shown below:



The figure shows the output layer of a neural network, labeled as  $o$ . The logits  $z_o^{(k)}$ , for  $k \in \{1,2,3\}$ , are the values in light green. These represent the outputs of the units before the activation function is applied. The vector  $\mathbf{z}$  is expressed as  $\mathbf{z}_o = [z_o^{(1)}, z_o^{(2)}, z_o^{(3)}]^\top$ .

For instance, the softmax for unit  $o,2$  in the figure is calculated as:

---

<sup>4</sup> Class labels in classification are arbitrary and unordered. You can assign numbers to the classes in any way, and the model's performance won't change as long as the mapping is consistent for all examples.

$$\text{softmax}(\mathbf{z}_o, 2) = \frac{e^{z_o^{(2)}}}{e^{z_o^{(1)}} + e^{z_o^{(2)}} + e^{z_o^{(3)}}}$$

The softmax function transforms a vector into a **discrete probability distribution (DPD)**, ensuring that  $\sum_{k=1}^D \text{softmax}(\mathbf{z}, k) = 1$ . A DPD assigns probabilities to values in a finite set, with their sum equaling 1. A **finite set** contains a countable number of distinct elements. For instance, in a classification task with classes 1, 2, and 3, these classes constitute a finite set. The softmax function maps each class to a probability, with these probabilities summing to 1.

Let's compute the probabilities step by step. Assume we have three logits,  $\mathbf{z} = [2.0, 1.0, 0.5]^\top$ , representing a document's classification into cinema, music, or science.

First, calculate  $e^{z^{(k)}}$  for each logit:

$$\begin{aligned} e^{z^{(1)}} &= e^{2.0} &\approx 7.39, \\ e^{z^{(2)}} &= e^{1.0} &\approx 2.72, \\ e^{z^{(3)}} &= e^{0.5} &\approx 1.65 \end{aligned}$$

Next, sum these values:  $\sum_{j=1}^3 e^{z^{(j)}} = 7.39 + 2.72 + 1.65 \approx 11.76$ .

Now use the softmax formula,  $\text{softmax}(\mathbf{z}, k) = \frac{e^{z^{(k)}}}{\sum_{j=1}^3 e^{z^{(j)}}}$ , to compute the probabilities:

$$\Pr(\text{cinema}) = \frac{7.39}{11.76} \approx 0.63, \quad \Pr(\text{music}) = \frac{2.72}{11.76} \approx 0.23, \quad \Pr(\text{science}) = \frac{1.65}{11.76} \approx 0.14$$

The softmax outputs are better characterized as “probability scores” rather than true probabilities. While these values sum to one and resemble class likelihoods, they don’t represent statistical probabilities of classes. Neural networks tend to produce outputs that diverge from true class probabilities. By contrast, models like logistic regression and Naïve Bayes generate genuine probabilities. For simplicity, though, I will refer to these outputs as “probabilities” throughout this book.

The **cross-entropy** loss measures how well predicted probabilities match the true distribution. The true distribution is typically a **one-hot vector** with a single element equal to 1 (the correct class) and zeros elsewhere. For example, a one-hot encoding with three classes looks like:

Class	One-hot vector
1	[1,0,0] <sup>⊤</sup>
2	[0,1,0] <sup>⊤</sup>
3	[0,0,1] <sup>⊤</sup>

The cross-entropy loss for a single example is:

$$\text{loss}(\tilde{\mathbf{y}}, \mathbf{y}) = - \sum_{k=1}^C y^{(k)} \log(\tilde{y}^{(k)}),$$

where  $C$  is the number of classes,  $\mathbf{y}$  is the one-hot encoded true label, and  $\tilde{\mathbf{y}}$  is the predicted probabilities. Here,  $y^{(k)}$  and  $\tilde{y}^{(k)}$  represent the  $k^{\text{th}}$  elements of  $\mathbf{y}$  and  $\tilde{\mathbf{y}}$ , respectively.

Since  $\mathbf{y}$  is one-hot encoded, only the term corresponding to the correct class contributes to the summation. The summation thus simplifies by retaining only that single term. Let's simplify it.

Suppose the correct class is  $c$ , so  $y^{(c)} = 1$  and  $y^{(k)} = 0$  for all  $k \neq c$ . In the summation, only the term where  $k = c$  will be non-zero. The equation simplifies to:

$$\text{loss}(\tilde{\mathbf{y}}, \mathbf{y}) = -\log(\tilde{y}^{(c)}) \quad (2.1)$$

This simplified form indicates that the loss corresponds to the negative logarithm of the probability assigned to the correct class. For  $N$  examples, the average loss is:

$$\text{loss} = -\frac{1}{N} \sum_{i=1}^N \log(\tilde{y}_i^{(c_i)}),$$

where  $c_i$  is the correct class index for the  $i^{\text{th}}$  example.

When used with softmax in the output layer, cross-entropy loss guides the network to assign high probabilities to correct classes while reducing probabilities for incorrect ones.

For a document classification example with three classes (cinema, music, and science), the network generates three logits. These logits are passed through the softmax function to convert them into probabilities for each class. The cross-entropy loss is then calculated between these scores and the true one-hot encoded labels.

Let's illustrate this by training a simple 2-layer neural network to classify documents into three classes. We begin by importing dependencies, setting a random seed, and defining the dataset:

```
import re, torch, torch.nn as nn

torch.manual_seed(42) ❶

docs = [
    "Movies are fun for everyone.",
    "Watching movies is great fun.",
    ...
    "Listen to folk music!"
]

labels = [1, 1, 1, 3, 3, 3, 2, 2, 2, 2]
num_classes = len(set(labels))
```

Setting the random seed in line ❶ ensures consistent random number generation across PyTorch runs. This guarantees **reproducibility**, allowing you to attribute performance changes to code or hyperparameter modifications rather than random variations. Reproducibility is also essential for teamwork, enabling collaborators to examine issues under identical conditions.

Next, we convert documents into a bag of words using two methods: `tokenize`, which splits input text into lowercase words, and `get_vocabulary`, which constructs the vocabulary:

```
def tokenize(text):
    return re.findall(r"\w+", text.lower()) ❶

def get_vocabulary(texts):
    tokens = {token for text in texts for token in tokenize(text)} ❷
    return {word: idx for idx, word in enumerate(sorted(tokens))} ❸

vocabulary = get_vocabulary(docs)
```

In line ❶, the regular expression `\w+` extracts individual words from the text. A **regular expression** is a sequence of characters used to define a search pattern. The pattern `\w+` matches sequences of “word characters,” such as letters, digits, and underscores.

The `findall` function from Python’s `re` module applies the regular expression and returns a list of all matches in the input string. In this case, it extracts all words.

In line ❷, the corpus is converted into a set of tokens by iterating through each document and extracting words using the same regular expression. In line ❸, these tokens are sorted alphabetically and mapped to unique indices, forming a vocabulary.

Once the vocabulary is built, the next step is to define the feature extraction function that converts a document into a feature vector:

```
def doc_to_bow(doc, vocabulary):
    tokens = set(tokenize(doc))
    bow = [0] * len(vocabulary)
    for token in tokens:
        if token in vocabulary:
            bow[vocabulary[token]] = 1
    return bow
```

The `doc_to_bow` function takes a document string and a vocabulary and returns the bag-of-words representation of the document.

Now, let’s transform our documents and labels into numbers:

```
vectors = torch.tensor(
    [doc_to_bow(doc, vocabulary) for doc in documents],
    dtype=torch.float32
)
labels = torch.tensor(labels, dtype=torch.long) - 1 ❶
```

Read first, buy later

42

The `vectors` tensor with shape  $(10, 26)$  represents 10 documents as rows and 26 vocabulary tokens as columns, while the `labels` tensor of shape  $(10,)$  contains the class label for each document. The labels use integer indices rather than one-hot encoding since PyTorch's cross-entropy loss function (`nn.CrossEntropyLoss`) expects this format.

Line ❶ uses `torch.long` to cast labels to 64-bit integers. The `-1` adjustment converts our original classes 1, 2, 3 to indices 0, 1, 2, which aligns with PyTorch's expectation that class indices begin at 0 for models and loss functions like `CrossEntropyLoss`.

PyTorch provides two APIs for model definition: the **sequential API** and the **module API**. While we used the straightforward `nn.Sequential` API to define our model in Section 1.8, we'll now explore building a multilayer perceptron using the more versatile `nn.Module` API:

```
input_dim = len(vocabulary)
hidden_dim = 50
output_dim = num_classes

class SimpleClassifier(nn.Module):
    def __init__(self, input_dim, hidden_dim, output_dim):
        super().__init__()
        self.fc1 = nn.Linear(input_dim, hidden_dim)
        self.relu = nn.ReLU()
        self.fc2 = nn.Linear(hidden_dim, output_dim)

    def forward(self, x):
        x = self.fc1(x) ❶
        x = self.relu(x) ❷
        x = self.fc2(x) ❸
        return x

model = SimpleClassifier(input_dim, hidden_dim, output_dim)
```

The `SimpleClassifier` class implements a **feedforward neural network** with two layers. Its constructor defines the network components:

1. A fully connected layer, `self.fc1`, maps the input of size `input_dim` (equal to the vocabulary size) to 50 (`hidden_dim`) outputs.
2. A ReLU activation function to introduce non-linearity.
3. A second fully connected layer, `self.fc2`, reduces the 50 intermediate outputs to `output_dim`, the number of unique labels.

The `forward` method describes the **forward pass**, where inputs flow through the layers:

- In line ❶, the input `x` of shape  $(10, 26)$  is passed to the first fully connected layer, transforming it to shape  $(10, 50)$ .
- In line ❷, output from this layer is fed through the ReLU activation function, keeping the shape  $(10, 50)$ .

- In line ❸, the result is sent to the second fully connected layer, reducing it from shape  $(10, 50)$  to  $(10, 3)$ , producing the model's final output with logits.

The `forward` method is called automatically when you pass input data to the model instance, like this: `model(input)`.

While `SimpleClassifier` omits a final softmax layer, this is intentional—PyTorch's `CrossEntropyLoss` combines softmax and cross-entropy loss internally for better numerical stability. This design eliminates the need for an explicit softmax in the model's forward pass.

With our model defined, the next steps, as outlined in Section 1.8, are to define the loss function, choose the gradient descent algorithm, and set up the training loop:

```
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.SGD(model.parameters(), lr=0.001)

for step in range(3000):
    optimizer.zero_grad()
    loss = criterion(model(vectors), labels)
    loss.backward()
    optimizer.step()
```

As you can see, the training loop is identical to the one in Section 1.8. Once the training is complete, we can test the model on a new document:

```
new_docs = [
    "Listening to rock music is fun.",
    "I love science very much."
]
class_names = ["Cinema", "Music", "Science"]

new_doc_vectors = torch.tensor(
    [doc_to_bow(new_doc, vocabulary) for new_doc in new_docs],
    dtype=torch.float32
)

with torch.no_grad(): ❶
    outputs = model(new_doc_vectors) ❷
    predicted_ids = torch.argmax(outputs, dim=1) + 1 ❸

for i, new_doc in enumerate(new_docs):
    print(f'{new_doc}: {class_names[predicted_ids[i].item() - 1]}')
```

Output:

```
Listening to rock is fun.: Music
I love scientific research.: Science
```

The `torch.no_grad()` statement in line ❶ disables the default gradient tracking. While gradients are essential during **training** to update model parameters, they're unnecessary during **testing** or **inference**. Since these phases don't involve parameter updates, disabling gradient tracking conserves memory and speeds up computation. Note that the terms "inference" and "evaluation" are often used interchangeably when referring to generating predictions on unseen data.

In line ❷, the model processes all inputs simultaneously during inference, just as it does during training. This parallel processing approach leverages vectorized operations, substantially reducing computation time compared to processing inputs one by one.

We only care about the final label, not the logits returned by the model. In line ❸, `torch.argmax` identifies the highest logit's index, corresponding to the predicted class. Adding 1 compensates for the earlier shift from 1-based to 0-based indexing.

While the bag-of-words approach offers simplicity and practicality, it has notable limitations. Most significantly, it fails to capture token order or context. Consider how "the cat chased the dog" and "the dog chased the cat" yield identical representations, despite conveying opposite meanings.

**N-grams** provide one solution to this challenge. An n-gram consists of  $n$  consecutive tokens from text. Consider the sentence "Movies are fun for everyone"—its bigrams (2-grams) include "Movies are," "are fun," "fun for," and "for everyone." By preserving sequences of tokens, n-grams retain contextual information that individual tokens cannot capture.

However, using n-grams comes at a cost. The vocabulary expands considerably, increasing the computational cost of model training. Additionally, the model requires larger datasets to effectively learn weights for the expanded set of possible n-grams.

Another limitation of bag-of-words is how it handles out-of-vocabulary words. When a word appears during inference that wasn't present during training—and thus isn't in the vocabulary—it can't be represented in the feature vector. Similarly, the approach struggles with synonyms and near-synonyms. Words like "movie" and "film" are processed as completely distinct terms, forcing the model to learn separate parameters for each. Since labeled data is often costly to obtain, resulting in small labeled datasets, it would be more efficient if the model could recognize and collectively process words with similar meanings.

**Word embeddings** address this challenge by representing semantically similar words as similar vectors.

## 2.2. Word Embeddings

Consider document 3 ("Enjoy a great movie today.") from earlier. We can break down this bag of words (BoW) into one-hot vectors representing individual words:

BoW	1	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	0	0	0	0	0
enjoy	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

<b>BoW</b>	1	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0	0
a	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
great	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
movie	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0
today	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0

As you can see, a bag-of-word vector of a document is a sum of one-hot vectors of its words. Now, let's examine the one-hot vectors and the BoW vector for the text "Films are my passion.":

<b>BoW</b>	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
films	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
are	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
my	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
passio	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
n																							

There are two key problems here. First, even when a word exists in the training data and vocabulary, one-hot encoding reduces it to a single 1 in a vector of zeros, giving the classifier almost no meaningful information to learn from.

Second, in the above document, most one-hot encoded word vectors add no value since three out of four become **zero vectors**—representing words missing from the vocabulary.

A better approach would let the model understand that "films," though unseen in training, shares semantic meaning with "movies." This would allow the feature vector for "films" to be processed similarly to "movies." Such an approach requires word representations that capture semantic relationships between words.

**Word embeddings** overcome the limitations of the bag-of-words model by representing words as **dense vectors** rather than **sparse** one-hot vectors. These lower-dimensional representations contain mostly non-zero values, with similar words having embeddings that exhibit high **cosine similarity**. The embeddings are learned from vast unlabeled datasets spanning millions to hundreds of millions of documents.

**Word2vec**, a widely-used embedding learning algorithm, exists in two variants. We'll examine the skip-gram formulation.

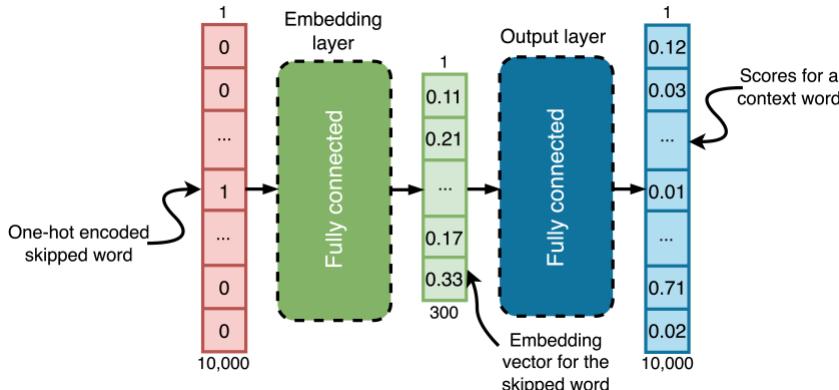
**Skip-grams** are formed by breaking text into sequences where one word is missing. For instance, in the sentence "*Professor Alan Turing's \* advanced computer science*," the skipped word, marked as \*, could be guessed as "work," "research," "contributions," "studies," or "theories." These words aren't exact synonyms but fit the context and share related meanings. A model can be trained to predict the missing word from its surrounding context words, and through this process, it learns the semantic relationships between words.

The process can also work in reverse: the skipped word can be used to predict its context words. This is the basis of the **skip-gram algorithm**.

The skip-gram size specifies how many context words are included. For a size of five, this means two words before and two after the skipped word. Here are examples of skip-grams of size five from our sentence, with different words skipped (marked as \*):

Skip-gram	Skipped word
professor alan * research advanced	turing's
alan turing's * advanced computer	research
turing's research * computer science	advanced

If the corpus vocabulary contains 10,000 words, the skip-gram model with an embedding layer of 300 units, is depicted below:



This is a skip-gram model with a skip-gram size of 5 and the embedding layer of 300 units. As you can see, the model uses a one-hot encoded skipped word to predict a context word, processing the input through two consecutive fully connected layers. It doesn't predict all context words at once but makes separate predictions for each.

Here's how it works for the skip-gram '*professor alan \* research advanced*' and the skipped word "turing's". We transform the skip-gram into 4 training pairs:

Skipped word (input)	Context word (target)	Position
turing's	professor	-2
turing's	alan	-1
turing's	research	+1
turing's	advanced	+2

For each pair of skipped and context words, say (turing's, professor), the model:

1. Takes "turing's" as input;
2. Converts it to a one-hot vector;
3. Passes it through the embedding layer to get the word embedding;
4. Passes the word embedding through the output layer;

## 5. Outputs probabilities for “professor.”

For a given context word, the output layer produces a probability vector across the vocabulary. Each value represents how likely that vocabulary word is to be the context word.

A curious reader might notice: if the input for each training pair remains constant—say, “turing’s”—why would the output differ? That’s a great observation! The output will indeed be identical for the same input. However, the loss calculation varies depending on each context word.

When using chat LMs, you may notice that the same question often yields different answers. While this might suggest the model is non-deterministic, that’s not accurate. An LLM is fundamentally a neural network, similar to a skip-gram model but with far more parameters. The apparent randomness comes from the way these models are *used* to generate text. During generation, words are sampled based on their predicted probabilities. Though higher-probability words are more likely to be chosen, lower-probability ones may still be selected. This sampling process creates the variations we observe in responses. We will talk about sampling in Chapter 5.

The skip-gram model uses **cross-entropy** as its loss function, just as in the three-class text classifier discussed earlier, but handles 10,000 classes—one for each word in the vocabulary. For each skip-gram in the training set, the model computes losses separately for each context word, such as the four words surrounding “turing’s,” then averages these losses to receive feedback on all context word predictions simultaneously.

This training approach enables the model to capture meaningful word relationships, even when working with the same input across different training pairs.

Here’s an example. For the input word “turing’s,” suppose the model assigns these probabilities to different vocabulary words: professor (0.1), alan (0.15), research (0.2), advanced (0.05). When training the model, each input-target word pair contributes to the loss function. For example, when “turing’s” appears with “professor” in the training data, the loss works to increase the score of 0.1. Similarly, when paired with “alan,” the loss works to increase 0.15, with “research” to increase 0.2, and with “advanced” to increase 0.05.

During backpropagation, the model adjusts its weights to make these scores higher for the given context words. For instance, the updated scores might be: professor: 0.11, alan: 0.17, research: 0.22, advanced: 0.07, while the scores for other vocabulary words decrease slightly.

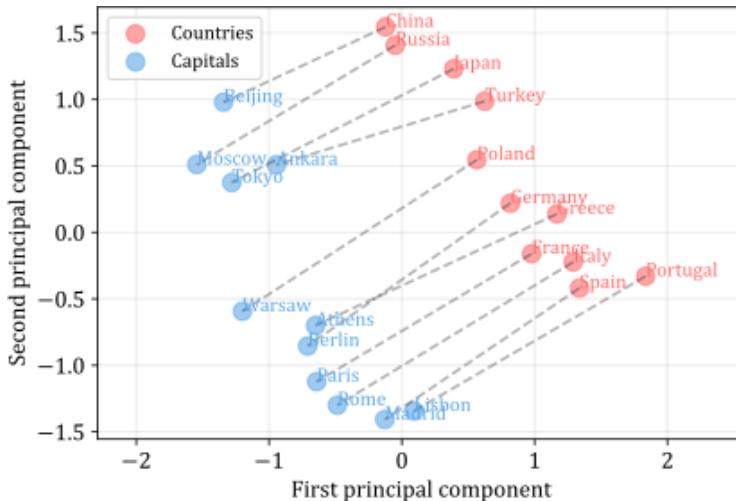
Once training is complete, the output layer is discarded. The embedding layer then serves as the new output layer. When given a one-hot encoded input word, the model produces a 300-dimensional vector—this is the word embedding.

Word2vec is just one method for learning word embeddings from large, unlabeled text corpora. Other methods, such as **GloVe** and **FastText**, offer alternative approaches, focusing on capturing global co-occurrence statistics or subword information to create more robust embeddings.

Using word embeddings to represent texts offers clear advantages over bag of words. One advantage is **dimensionality reduction**, which compresses the word representation from the size

of the vocabulary (as in one-hot encoding) to a small vector, typically between 100 and 1000 dimensions. This makes it feasible to process very large corpora in machine learning tasks.

**Semantic similarity** is another advantage of word embeddings. Words with similar meanings are mapped to vectors that are close to each other in the embedding space. For example, consider word embeddings trained by Google on a news corpus containing about 100 billion words.<sup>5</sup> In the graph below, “Moscow” and “Beijing,” or “Russia” and “China,” are represented by points located near one another. This reflects their semantic relationships:



The graph displays a two-dimensional projection of 300-dimensional word2vec embedding vectors for countries and their capitals. Words with related meanings cluster together, while nearly parallel lines connect cities to their respective countries, revealing their semantic relationships.

The skip-gram model captures semantic similarity when words occur in similar contexts, even without direct co-occurrence. For instance, if the model produces different probabilities for “films” and “movies,” the loss function drives it to predict similar ones, since context words for these terms frequently overlap. Through backpropagation, the embedding layer outputs for these words converge.

Before word embeddings, **WordNet** (created at Princeton in 1985) attempted to capture word relationships by organizing words into sets of synonyms and recording semantic links between them. While effective, these hand-crafted mappings couldn’t scale to large vocabularies or capture the subtle patterns in word usage that naturally emerge from embedding-based approaches.

<sup>5</sup> These embeddings can be found online by using the “GoogleNews-vectors-negative300.bin.gz” query. A backup is available on the book’s wiki at [thelmbbook.com/data/word-vectors](http://thelmbbook.com/data/word-vectors).

Because directly visualizing 300-dimensional vectors isn't possible, we used a **dimensionality reduction** technique called **principal component analysis (PCA)** to project them onto two dimensions, known as first and second **principal components**.

Dimensionality reduction algorithms compress high-dimensional vectors while maintaining their relationships. The first and second principal components in the above graph preserved the semantic connections between words, revealing their relationships.

For resources on PCA and other dimensionality reduction methods, check the recommended material listed on the book's wiki.

Word embeddings capture the meaning of words and their relationships to other words. They are fundamental to many natural language processing (NLP) tasks. Neural language models, for example, encode documents as matrices of word embeddings. Each row corresponds to a word's embedding vector, and its position reflects the word's position in the document.

The discovery that word2vec embeddings support meaningful arithmetic operations (like "king – man + woman ≈ queen") was a pivotal moment, revealing that neural networks could encode semantic relationships in a mathematical space where vector operations produced changes in word meaning. This made the invention of neural networks capable of doing complex math on words, like large language models do, only a matter of time.

Modern language models, though, often use subwords—tokens smaller than complete words. Before moving on to language models—the main topic of this book—let's first examine byte-pair encoding, a widely used subword tokenization method.

## 2.3. Byte-Pair Encoding

**Byte-pair encoding (BPE)** is a tokenization algorithm that addresses the challenges of handling out-of-vocabulary words by breaking words into smaller units called **subwords**.

Initially a data compression technique, BPE was adapted for NLP by treating words as sequences of characters. It merges the most frequent symbol pairs—characters or subwords—into new subword units. This continues until the vocabulary reaches the target size.

Below is the basic BPE algorithm:

1. **Initialization:** Use a text corpus. Split each word in the corpus into individual characters. For example, the word "hello" becomes "h e l l o". The initial vocabulary consists of all unique characters in the corpus.
2. **Iterative merging:**

- **Count adjacent symbol pairs:** Treat each character as a **symbol**. Go through the corpus and count every pair of adjacent symbols. For example, in “h e l l o”, the pairs are “h e”, “e l”, “l l”, “l o”.
- **Select the most frequent symbol pair:** Identify the symbol pair with the highest count across the entire corpus. For instance, if “l l” occurs most frequently, it will be selected.
- **Merge the selected pair:** Replace all occurrences of the most frequent symbol pair with a new single **merged symbol**. For example, “l l” would be replaced with a new merged symbol “ll”. The word “h e l l o” now becomes “h e ll o”.
- **Update the vocabulary:** Add the new merged symbol to the vocabulary. The vocabulary now includes the original characters and the new symbol “ll”.

3. **Repeat:** Continue the iterative merging until the vocabulary reaches the desired size.

The algorithm is simple, but implementing it directly on large corpora is inefficient. Recomputing symbol pairs or updating the entire corpus after each merge is computationally expensive.

A more efficient approach initializes the vocabulary with all unique words in the corpus and their counts. Pair counts are calculated using these word counts, and the vocabulary is updated iteratively by merging the most popular pairs. Let's write the code:

```
from collections import defaultdict

def initialize_vocabulary(corpus):
    vocabulary = defaultdict(int)
    charset = set()
    for word in corpus:
        word_with_marker = '_' + word ①
        characters = list(word_with_marker) ②
        charset.update(characters) ③
        tokenized_word = ' '.join(characters) ④
        vocabulary[tokenized_word] += 1 ⑤
    return vocabulary, charset
```

The function generates a vocabulary that represents words as sequences of characters and tracks their counts. Given a **corpus** (a list of words), it returns two outputs: **vocabulary**, a dictionary mapping each word—tokenized with spaces between characters—to its count, and **charset**, a set of all unique characters present in the corpus.

Here's how it works:

- Line ① adds a word boundary marker “\_” to the start of each word to differentiate subwords at the beginning from those in the middle. For example, “\_re” in “restart” is distinct from “re” in “agree.” This helps rebuild sentences from tokens generated using

the model. When a token starts with “\_”, it marks the beginning of a new word, requiring a space to be added before it.

- Line ❷ splits each word into individual **characters**.
- Line ❸ updates **charset** with any new characters encountered in the word.
- Line ❹ joins **characters** with spaces to create a tokenized version of the word. For example, the word “hello” becomes `_ h e l l o`.
- Line ❺ adds **tokenized\_word** to **vocabulary** with its count incremented.

After the initialization, BPE iteratively merges the most frequent pairs of tokens (bigrams) in the vocabulary. By removing spaces between these pairs, it forms progressively longer tokens.

```
def get_pair_counts(vocabulary):
    pair_counts = defaultdict(int)
    for tokenized_word, count in vocabulary.items():
        tokens = tokenized_word.split() ❶
        for i in range(len(tokens) - 1):
            pair = (tokens[i], tokens[i + 1]) ❷
            pair_counts[pair] += count ❸
    return pair_counts
```

The function counts how often adjacent token pairs appear in the tokenized vocabulary words. The input **vocabulary** maps tokenized words to their counts, and the output is a dictionary of token pairs and their total counts.

For each **tokenized\_word** in **vocabulary**, we split it into **tokens** in line ❶. A nested loop forms adjacent token pairs in line ❷ and increments their count by the word’s count in line ❸.

```
def merge_pair(vocabulary, pair):
    new_vocabulary = {}
    bigram = re.escape(' '.join(pair)) ❶
    pattern = re.compile(r"(?<!\S)" + bigram + r"(?!\\S)") ❷
    for tokenized_word, count in vocabulary.items():
        new_tokenized_word = pattern.sub("".join(pair), tokenized_word) ❸
        new_vocabulary[new_tokenized_word] = count
    return new_vocabulary
```

The function merges the input token pair in all tokenized words from the vocabulary. It returns a new vocabulary where every occurrence of the **pair** is merged into a single token. For example, if the pair is `('e', 'l')` and a tokenized word is `"_ h e l l o"`, merging ‘e’ and ‘l’ removes the space between them, resulting in `"_ h el l o"`.

In line ❶, the `re.escape` function automatically adds backslashes to special characters in a string (like `.`, `*`, or `?`), so they are interpreted as literal characters rather than having their special meaning in regular expressions.

The regular expression in line ❷ matches only whole token pairs. It ensures the **bigram** is not part of a larger word by checking for the absence of non-whitespace characters immediately before and

after the match. For instance "good morning" matches in "this is good morning", but not in "thisisgood morning", where "good" is part of "thisisgood".

The expression `(?<!\\S)` is a **negative lookbehind** assertion. It checks that no non-whitespace character (`\\S`) comes before the bigram. This ensures the bigram is either at the start of the `tokenized_word` or follows a whitespace character (like a space, tab, or newline), preventing it from being part of a longer word on the left. Similarly, `(?!\\S)` is a **negative lookahead**. It ensures no non-whitespace character follows the bigram, meaning the bigram is followed by whitespace, punctuation, or the end of the line, ensuring it is not part of a longer word on the right.

Finally, in line ❸, the function uses `pattern.sub()` to replace all occurrences of the matched pattern with the joined pair, creating the new tokenized word.

The function below implements the BPE algorithm, merging the most frequent token pairs iteratively until no merges remain or the target vocabulary size is reached:

```
def byte_pair_encoding(corpus, vocab_size):
    vocabulary, charset = initialize_vocabulary(corpus)
    merges = []
    tokens = set(charset)
    while len(tokens) < vocab_size: ❶
        pair_counts = get_pair_counts(vocabulary)
        if not pair_counts: ❷
            break
        most_frequent_pair = max(pair_counts, key=pair_counts.get) ❸
        merges.append(most_frequent_pair)
        vocabulary = merge_pair(vocabulary, most_frequent_pair) ❹
        new_token = ''.join(most_frequent_pair) ❺
        tokens.add(new_token) ❻

    return vocabulary, merges, charset, tokens
```

This function processes a corpus to produce the components needed for a tokenizer. It initializes the vocabulary and character set, creates an empty `merges` list for storing merge operations, and sets `tokens` to the initial character set. Over time, `tokens` grows to include all unique tokens the tokenizer will be able to generate.

The loop in line ❶ continues until the number of tokens supported by the tokenizer reaches `vocab_size` or no pairs remain to merge. Line ❷ checks if there are no more valid pairs, in which case the loop exits. Line ❸ finds the most frequent token pair, which is merged throughout the vocabulary in line ❹ to create a new token in line ❺. This new token is added to the `tokens` set in line ❻, and the merge is recorded in `merges`.

The function returns four outputs: the updated vocabulary, the list of merge operations, the original character set, and the final set of unique tokens.

The function below tokenizes a word using a trained tokenizer:

```
def tokenize_word(word, merges, vocabulary, charset, unk_token=<UNK>):
    word = ' ' + word
    if word in vocabulary:
        return [word]
    tokens = [char if char in charset else unk_token for char in word]

    for left, right in merges:
        i = 0
        while i < len(tokens) - 1:
            if tokens[i:i+2] == [left, right]:
                tokens[i:i+2] = [left + right]
            else:
                i += 1
    return tokens
```

This function tokenizes a word using `merges`, `vocabulary`, and `charset` from `byte_pair_encoding`. The word is first prefixed. If the prefixed word exists in the `vocabulary`, it returns it as the only token. Otherwise, the word is split into characters, with any character not in `charset` replaced by `unk_token`. These characters are then iteratively merged using the order of rules in `merges`.

To tokenize a text, we first split it into words based on spaces and then tokenize each word individually. The [thelmbbook.com/nb/2.1](http://thelmbbook.com/nb/2.1) notebook contains code for training a BPE tokenizer by using a news corpus.

The tokenized version of the sentence “*Let’s proceed to the language modeling chapter.*” using the tokenizer trained in the notebook, is:

```
["_Let", "", "s", "_proceed", "_to", "_the", "_language", "_model", "ing",
"_part", "."]
```

Here, “let’s,” and “modeling,” were broken into subwords. This indicates their relative rarity in the training data and a small target vocabulary size of 5000 tokens.

As you might have observed, the `tokenize_word` algorithm is inefficient, because it iterates over all merges in line ④ and checks every token pair in line ⑤. This creates slow nested loops. However, with vocabularies often exceeding 100,000 tokens in modern language models, most input words are already in the vocabulary, avoiding subword tokenization altogether. The notebook includes an optimized version of `tokenize_word` that uses caching and a precomputed data structure to eliminate nested loops. Running this optimized version on the same input reduces tokenization time from 0.0549 seconds to 0.0037 seconds. While the exact numbers depend on your system and input, the optimized approach will consistently be faster.

For languages without spaces, like Chinese, or for multilingual models, the initial space-based tokenization is typically skipped. Instead, the text is split into individual characters. From there, BPE proceeds as usual, merging the most frequent character or token pairs to form subwords.

We're now ready to explore the core ideas and methods of language modeling. We'll start with traditional count-based methods and cover neural network-based techniques in later chapters.

## 2.4. Language Model

A **language model** predicts the next token in a sequence by estimating its conditional probability based on previous tokens. It assigns a probability to all possible next tokens, enabling the selection of the most likely one. This capability supports tasks like text generation, machine translation, and speech recognition. Trained on large unlabeled text corpora, language models learn statistical patterns in language, allowing them to generate human-like text.

Formally, for a sequence  $\mathbf{s}$  of  $L$  tokens  $(t_1, t_2, \dots, t_L)$ , a language model computes:

$$\Pr(t = t_{L+1} | \mathbf{s} = (t_1, t_2, \dots, t_L)) \quad (2.2)$$

Here,  $\Pr$  represents the conditional probability distribution over the vocabulary for the next token. A **conditional probability** quantifies the likelihood of one event occurring given that another has already occurred. In language models, it reflects the probability of a specific token being the next one, given the preceding sequence of tokens. This sequence is often referred to as the **input sequence, context, or prompt**.

The following notations are equivalent to Equation 2.2:

$$\Pr(t_{L+1} | t_1, t_2, \dots, t_L) \text{ or } \Pr(t_{L+1} | \mathbf{s}) \quad (2.3)$$

We will select different notations, ranging from concise to detailed, based on the context.

For any token  $t$  and sequence  $\mathbf{s}$ , the conditional probability satisfies  $\Pr(t|\mathbf{s}) \geq 0$ , meaning probabilities are always non-negative. Furthermore, the probabilities for all possible next tokens in the vocabulary  $\mathcal{V}$  must sum to 1:  $\sum_{t \in \mathcal{V}} P(t|\mathbf{s}) = 1$ . This ensures the model outputs a valid probability distribution over the vocabulary.

To illustrate, let's consider an example with a vocabulary  $\mathcal{V}$  containing 5 words: "are," "cool," "language," "models," and "useless." For the sequence  $\mathbf{s} = (\text{language, models, are})$ , a language model could output the following probabilities for each possible next word in  $\mathcal{V}$ :

$$\begin{aligned} \Pr(t = \text{are} | \mathbf{s} = (\text{language, models, are})) &= 0.01 \\ \Pr(t = \text{cool} | \mathbf{s} = (\text{language, models, are})) &= 0.77 \\ \Pr(t = \text{language} | \mathbf{s} = (\text{language, models, are})) &= 0.02 \\ \Pr(t = \text{models} | \mathbf{s} = (\text{language, models, are})) &= 0.15 \\ \Pr(t = \text{useless} | \mathbf{s} = (\text{language, models, are})) &= 0.05 \end{aligned}$$

The illustration demonstrates how the language model assigns probabilities across its vocabulary for each potential next word, with “cool” receiving the highest probability. These probabilities sum to 1, forming a valid discrete probability distribution.

This type of model is an **autoregressive language model**, also known as a **causal language model**. **Autoregression** involves predicting an element in a sequence using only its predecessors. Such models excel at text generation and include Transformer-based **chat LMs**, and all language models discussed in this book.

In contrast, **masked language models**, such as BERT—a pioneering Transformer-based model—use a different approach. These models predict intentionally masked tokens within sequences, utilizing both preceding and following context. This bidirectional approach particularly suits tasks like text classification and named entity recognition.

Before neural networks became standard for language modeling, traditional methods relied on statistical techniques. These count-based models, still used in tools like smartphone autocomplete systems, estimate the probability of word sequences based on word or n-gram frequency counts learned from a corpus. To understand these methods better, let’s implement a simple count-based language model.

## 2.5. Count-Based Language Model

We’ll focus on a trigram model ( $n = 3$ ) to illustrate how this works. In a trigram model, the probability of a token is calculated based on the two preceding tokens:

$$\Pr(t_i | t_{i-2}, t_{i-1}) = \frac{C(t_{i-2}, t_{i-1}, t_i)}{C(t_{i-2}, t_{i-1})}, \quad (2.4)$$

where  $C(\cdot)$  denotes the count of occurrences of an n-gram in the training data.

For instance, if the trigram “language models rock” appears 50 times in the corpus and “language models” appears 200 times overall, then:

$$\Pr(\text{rock} | \text{language, models}) = \frac{50}{200} = 0.25$$

This means that “rock” follows “language models” 25% of the time in our training data.

Equation 2.4 is the **maximum likelihood estimate (MLE)** of a token’s probability given its context. It calculates the relative frequency of a trigram compared to all trigrams with the same two-token history. As the training corpus grows, the MLE provides a better estimate for n-grams that appear frequently in the data. This reflects a basic statistical principle: larger datasets lead to more accurate estimates.

However, a limited-size corpus poses a problem: some n-grams we may encounter in practice might not appear in the training data. For instance, if the trigram “language models sing” never appears in our corpus, its probability would be zero according to the MLE:

$$\Pr(\text{sing}|\text{language, models}) = \frac{0}{200} = 0$$

This is problematic because it assigns a zero probability to any sequence containing an unseen n-gram, even if it is a valid phrase. To solve this, several techniques have been developed, one of which is **backoff**. The idea is simple: if a higher-order n-gram (e.g., trigram) is not observed, we “back off” to a lower-order n-gram (e.g., bigram). The probability  $\Pr(t_i|t_{i-2}, t_{i-1})$  is given by one of the following expressions, depending on whether the condition is true:

Expression	Condition
$\frac{C(t_{i-2}, t_{i-1}, t_i)}{C(t_{i-2}, t_{i-1})}$	if $C(t_{i-2}, t_{i-1}, t_i) > 0$
$\Pr(t_i t_{i-1})$	if $C(t_{i-2}, t_{i-1}, t_i) = 0$ and $C(t_{i-1}, t_i) > 0$
$\Pr(t_i)$	otherwise

Here,  $C(t_{i-2}, t_{i-1}, t_i)$  is the count of the trigram  $(t_{i-2}, t_{i-1}, t_i)$ ,  $C(t_{i-2}, t_{i-1})$  and  $C(t_{i-1}, t_i)$  are the counts of the bigrams  $(t_{i-2}, t_{i-1})$  and  $(t_{i-1}, t_i)$  respectively.

The bigram probability and unigram probability are computed as:

$$\Pr(t_i|t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}, \quad \Pr(t_i) = \frac{C(t_i) + 1}{W + V},$$

where  $C(t_i)$  is the count of the token  $t_i$ ,  $W$  is the total number of tokens in the corpus, and  $V$  is the vocabulary size.

Adding 1 to  $C(t_i)$ , known as **add-one smoothing** or **Laplace smoothing**, addresses zero probabilities for tokens not present in the corpus. If we used the actual frequency  $\Pr(t_i) = \frac{C(t_i)}{W}$ , any token not found in the corpus would have a zero probability, creating problems when the model encounters valid but unseen tokens.

Laplace smoothing solves this by adding 1 to each token count, ensuring all tokens, including unseen ones, receive a small, non-zero probability. The denominator is adjusted by adding  $V$  to account for the extra counts introduced in the numerator.

Now, let's implement a language model with backoff in the `CountLanguageModel` class:

```
class CountLanguageModel:
    def __init__(self, n): ❶
        self.n = n
        self.ngram_counts = [{} for _ in range(n)] ❷
        self.total_unigrams = 0

    def predict_next_token(self, context): ❸
        for n in range(self.n, 1, -1): ❹
```

```

    if len(context) >= n - 1: ❸
        context_n = tuple(context[-(n - 1):]) ❹
        counts = self.ngram_counts[n - 1].get(context_n)
        if counts:
            return max(counts.items(), key=lambda x: x[1])[0]
    unigram_counts = self.ngram_counts[0].get(())
    if unigram_counts:
        return max(unigram_counts.items(), key=lambda x: x[1])[0]
    return None

```

In line ❶, the model is initialized with an `n` parameter, defining the maximum n-gram order (e.g., `n=3` for trigrams). The `ngram_counts` list in line ❷ stores n-gram frequency dictionaries for unigrams, bigrams, trigrams, etc., populated during training. For `n=3`, given the corpus “*Language models are powerful. Language models are useful.*” in lowercase with punctuation removed, `self.ngram_counts` would contain:

```

ngram_counts[0] = {(): {"language": 2, "models": 2, "are": 2, "powerful": 1,
                       "useful": 1}}

ngram_counts[1] = {("language",): {"models": 2}, ("models",): {"are": 2},
                   ("are",): {"powerful": 1, "useful": 1}}

ngram_counts[2] = {("language", "models"): {"are": 2}, ("models", "are"): {"powerful": 1,
                           "useful": 1}}

```

The `predict_next_token` method uses backoff to predict the next token. Starting from the highest n-gram order in line ❷, it checks if the context contains enough tokens for this n-gram order in line ❸. If so, it extracts the relevant context in line ❹ and attempts to find a match in `ngram_counts`. If no match is found, it backs off to lower-order n-grams or defaults to unigram counts. For instance, given `context=["language", "models", "are"]` and `n=3`:

- First iteration: `context_n = ("models", "are")`
- Second iteration (if needed): `context_n = ("are",)`
- Last resort: unigram counts with empty tuple key ()

If a matching context is found, the method returns the token with the highest frequency for that context. For input `["language", "models"]` it will return `"are"`, the token with highest count among values for the key `("language", "models")` in `ngram_counts[2]`. However, for the input `["english", "language"]` it will not find the key `("english", "language")` in `ngram_counts[2]`, so it will backoff to `ngram_counts[1]` and return `"models"`, the token with highest count among values for the key `("language",)`.

Now, let's define the method that trains our model:

```

def train(model, tokens):
    model.total_unigrams = len(tokens)
    for n in range(1, model.n + 1): ❶

```

```

counts = model.ngram_counts[n - 1]
for i in range(len(tokens) - n + 1):
    context = tuple(tokens[i:i + n - 1]) ②
    next_token = tokens[i + n - 1] ③
    if context not in counts:
        counts[context] = defaultdict(int)
    counts[context][next_token] = counts[context][next_token] + 1

```

The `train` method takes a model (an instance of `CountLanguageModel`) and a list of tokens (the training corpus) as input. It populates the n-gram counts in the model using these tokens.

In line ①, the method iterates over n-gram orders from 1 to `model.n` (inclusive). For each `n`, it generates n-grams of that order from the token sequence and counts their occurrences.

Lines ② and ③ extract contexts and their subsequent tokens to build a nested dictionary where each context maps to a dictionary of following tokens and their counts. These counts are stored in `model.ngram_counts`, which the `predict_next_token` method later uses to make predictions based on context.

Now, let's train the model:

```

set_seed(42)
n = set_hyperparameters()
data_url = "https://www.thelmbbook.com/data/brown"
train_corpus, test_corpus = download_and_prepare_data(data_url)

model = CountLanguageModel(n)
train(model, train_corpus)

perplexity = compute_perplexity(model, test_corpus)
print(f"\nPerplexity on test corpus: {perplexity:.2f}")

contexts = [
    "i will build a",
    "the best place to",
    "she was riding a"
]

for context in contexts:
    words = tokenize(context)
    next_word = model.predict_next_token(words)
    print(f"\nContext: {context}")
    print(f"Next token: {next_word}")

```

The full implementation of this model, including the methods to retrieve and process the training data, can be found in the [thelmbbook.com/nb/2.2](https://www.thelmbbook.com/nb/2.2) notebook. Within the `download_and_prepare_data` method, the corpus is downloaded, converted to lowercase,

tokenized into words, and divided into **training** and **test** partitions with a 90/10 split. Let's take a moment to understand why this last step is critical.

In machine learning, using the entire dataset for training leaves no way to evaluate whether the model **generalizes** well. A frequent issue is **overfitting**, where the model excels on training data but struggles to make accurate predictions on unseen, new data.

Partitioning the dataset into training and test sets is a standard practice to control overfitting. It involves two steps: (1) shuffling the data and (2) splitting it into two subsets. The larger subset, called the training data, is used for training the model, while the smaller subset, called the test data, is used to evaluate the model's performance on unseen examples.

The test set requires sufficient size to reliably estimate model performance. A test ratio of 0.1 to 0.3 (10% to 30% of the entire dataset) is common, though this varies with dataset size. For very large datasets, even a smaller test set can provide reliable performance estimates.

The training data comes from the **Brown Corpus**, a collection of over 1 million words from American English texts published in 1961. This corpus is frequently used in linguistic studies.

When you run the code, you will see the following output:

```
Perplexity on test corpus: 302.08
```

```
Context: i will build a  
Next word: wall
```

```
Context: the best place to  
Next word: live
```

```
Context: she was riding a  
Next word: horse
```

Ignore the perplexity number for now; we'll discuss it shortly. Count-based language models can produce reasonable immediate continuations, making them effective for autocomplete systems. However, they have notable limitations. These models generally work with word-tokenized corpora, as their n-gram size is typically small (up to  $n = 5$ ). Extending beyond this would require too much memory and lead to slower processing. Subword tokenization, while potentially more efficient, often results in many n-grams that represent only fragments of words, degrading the quality of next-word predictions.

Word-level tokenization creates a significant drawback: count-based models cannot handle out-of-vocabulary (OOV) words. This is similar to the issue seen in the **bag-of-words** approach discussed in Section 2.1. For example, consider the context: "according to WHO, COVID-19 is a." If "COVID-19" wasn't in the training data, the model would back off repeatedly until it relies only on "is a." This severely limits the context and makes meaningful predictions unlikely.

Short contexts make it impossible for count-based models to capture long-range dependencies in language. While modern Transformer models also have context size limits, these can handle thousands of tokens. In contrast, training a count-based model with a long context, such as 1000 tokens, would require storing counts for all n-grams from  $n = 1$  to  $n = 1000$ , which is infeasible due to memory constraints.

Another limitation is that count-based models cannot be adapted for downstream tasks after training. Their n-gram counts are fixed, and any adjustment requires retraining on new data.

These limitations have led to the development of advanced methods, particularly neural network-based language models, which solve these problems and have largely replaced count-based models in modern natural language processing. Approaches like recurrent neural networks and transformers, which we will discuss in the next two chapters, handle longer contexts effectively, producing more coherent and context-aware text. Before exploring these methods, let's look at how to evaluate a language model's quality.

## 2.6. Evaluating Language Models

Evaluating language models measures their performance and allows comparing models. Several metrics and techniques are commonly used. Let's look at the main ones.

### 2.6.1. Perplexity

**Perplexity** is a widely used metric for evaluating language models. It measures how well a model predicts a text. Lower perplexity values indicate a better model—one that is more confident in its predictions. Perplexity is defined as the exponential of the average **negative log-likelihood** per token in the test set:

$$\text{Perplexity}(\mathcal{D}, k) = \exp\left(-\frac{1}{D} \sum_{i=1}^D \log \Pr(t_i | t_{\max(1, i-k)}, \dots, t_{i-1})\right) \quad (2.5)$$

Here,  $\mathcal{D}$  represents the test set,  $D$  is the total number of tokens in it,  $t_i$  is the  $i^{\text{th}}$  token, and  $\Pr(t_i | t_{\max(1, i-k)}, \dots, t_{i-1})$  is the probability the model assigns to  $t_i$  given its preceding context window of size  $k$ , where  $\max(1, i - k)$  ensures we start from the first token when there aren't enough preceding tokens to fill the context window. The notations  $\exp(x)$  and  $e^x$ , where  $e$  is **Euler's number**, are equivalent.

Perplexity can be interpreted as the number of choices a model considers when predicting the next token. For example, a perplexity of 10 means the model behaves as if choosing from 10 equally likely options on average.

If a language model assigns equal probability to every token in a vocabulary of size  $V$ , its perplexity equals  $V$ . This represents an intuitive upper bound—no model can be more uncertain than assigning equal likelihood to all words.

Let's calculate perplexity using the example text with word-level tokenization: "*We are evaluating a language model for English.*" To keep things simple, we assume a context of three words. We begin by determining the probability of each word based on its preceding context of three words as provided by the model. Here are the probabilities:

$$\begin{aligned}
 \Pr(\text{We}) &= 0.10 \\
 \Pr(\text{are} \mid \text{We}) &= 0.20 \\
 \Pr(\text{evaluating} \mid \text{We, are}) &= 0.05 \\
 \Pr(\text{a} \mid \text{We, are, evaluating}) &= 0.50 \\
 \Pr(\text{language} \mid \text{are, evaluating, a}) &= 0.30 \\
 \Pr(\text{model} \mid \text{evaluating, a, language}) &= 0.40 \\
 \Pr(\text{for} \mid \text{a, language, model}) &= 0.15 \\
 \Pr(\text{English} \mid \text{language, model, for}) &= 0.25
 \end{aligned}$$

Using the probabilities, we compute the negative log-likelihood for each word:

$$\begin{aligned}
 -\log(P(\text{We})) &= -\log(0.10) \approx 2.30 \\
 -\log(P(\text{are} \mid \text{We})) &= -\log(0.20) \approx 1.61 \\
 -\log(P(\text{evaluating} \mid \text{We, are})) &= -\log(0.05) \approx 3.00 \\
 -\log(P(\text{a} \mid \text{We, are, evaluating})) &= -\log(0.50) \approx 0.69 \\
 -\log(P(\text{language} \mid \text{are, evaluating, a})) &= -\log(0.30) \approx 1.20 \\
 -\log(P(\text{model} \mid \text{evaluating, a, language})) &= -\log(0.40) \approx 0.92 \\
 -\log(P(\text{for} \mid \text{a, language, model})) &= -\log(0.15) \approx 1.90 \\
 -\log(P(\text{English} \mid \text{language, model, for})) &= -\log(0.25) \approx 1.39
 \end{aligned}$$

Next, we sum these values and divide the sum by the number of words (8) to get the average:

$$(2.30 + 1.61 + 3.00 + 0.69 + 1.20 + 0.92 + 1.90 + 1.39)/8 \approx 1.63$$

Finally, we exponentiate the average negative log-likelihood to obtain the perplexity:

$$e^{1.63} \approx 5.10$$

So, the model's perplexity on this text, using a 3-word context, is about 5.10. This means that, on average, the model acts as if it selects from roughly 5 equally likely options for each word in the given context. A perplexity of 5.10 is relatively low, showing the model is confident in its predictions for this sequence. But for a proper evaluation, perplexity must be measured on a much larger test set.

Now, let's calculate the perplexity of the count-based model from the previous section. To do this, the model must be updated to return the probability of a token given a specific context. Add this function to the `CountLanguageModel` we implemented earlier:

```
def get_probability(self, token, context):
    for n in range(self.n, 1, -1): ❶
        if len(context) >= n - 1:
            context_n = tuple(context[-(n - 1):])
```

```

counts = self.ngram_counts[n - 1].get(context_n)
if counts: ❷
    total = sum(counts.values()) ❸
    count = counts.get(token, 0)
    if count > 0:
        return count / total ❹
unigram_counts = self.ngram_counts[0].get() ❺
count = unigram_counts.get(token, 0)
V = len(unigram_counts)
return (count + 1) / (self.total_unigrams + V) ❻

```

The `get_probability` function is similar to `predict_next_token`. Both loop through n-gram orders in reverse (line ❶) and extract the relevant context (`context_n`). If `context_n` matches in the n-gram counts (line ❷), the function retrieves the token counts. If no match exists, it backs off to lower-order n-grams and, finally, unigrams (line ❺).

Unlike `predict_next_token`, which finds the most probable token, `get_probability` calculates a token's probability. In line ❸, `total` is the sum of counts for tokens following the context, acting as the denominator. Line ❹ divides the token count by `total` to compute its probability. If no higher-order match exists, line ❻ uses add-one smoothing with unigram counts.

The `calculate_perplexity` method computes a language model's perplexity for a token sequence. It takes three arguments: the model, the token sequence, and the context size:

```

def compute_perplexity(model, tokens, context_size):
    if not tokens:
        return float('inf')
    total_log_likelihood = 0
    num_tokens = len(tokens)
    for i in range(num_tokens): ❶
        context_start = max(0, i - context_size)
        context = tuple(tokens[context_start:i]) ❷
        word = tokens[i]
        probability = model.get_probability(word, context)
        total_log_likelihood += math.log(probability) ❸
    average_log_likelihood = total_log_likelihood / num_tokens ❹
    perplexity = math.exp(-average_log_likelihood) ❺
    return perplexity

```

In line ❶, the function iterates through each token in the sequence. For every token:

- Line ❷ extracts its context, using up to `context_size` tokens before it. The expression `max(0, i - context_size)` ensures indices stay within bounds.
- In line ❸, the log of the token's probability is added to the cumulative log-likelihood. The `get_probability` method from the model handles the probability calculation.

Once all tokens are processed, line ④ computes the average log-likelihood by dividing the total log-likelihood by the number of tokens.

Finally, in line ⑤, the perplexity is computed as the exponential of the negative average log-likelihood, as described in Equation 2.5.

By applying this method to the `test_corpus` sequence from the [thelmbook.com/nb/2.2](http://thelmbook.com/nb/2.2) notebook, we observe the following output:

```
Perplexity on test corpus: 302.08
```

This perplexity is very high. For example, **GPT-2** has a perplexity of about 20, while modern LLMs achieve values below 10. Later, we'll compute perplexities for RNN and Transformer-based models and compare them with the perplexity of this count-based model.

## 2.6.2. ROUGE

Perplexity is a standard metric used to evaluate language models trained on large, unlabeled datasets by measuring how well they predict the next token in a given context. These models are referred to as **pretrained models** or **base models**. As we will discuss in the chapter on large language models, their ability to perform specific tasks or answer questions comes from **supervised finetuning**. This additional training uses a smaller, labeled dataset where input contexts are matched with target outputs, such as answers or task-specific results. This enables capabilities like translation or summarization.

Perplexity is not an ideal metric for evaluating a finetuned model. Instead, metrics are needed that compare the model's output to reference texts, often called **ground truth**. A common choice is **ROUGE** (Recall-Oriented Understudy for Gisting Evaluation). ROUGE is widely used for tasks like summarization and machine translation. It evaluates text quality by measuring overlaps, such as tokens or n-grams, between the generated text and the reference.

ROUGE has several variants, each focusing on different aspects of text similarity. Here, we'll discuss three widely used ones: ROUGE-1, ROUGE-N, and ROUGE-L.

**ROUGE-N** evaluates the overlap of n-grams between the generated and reference texts, with N indicating the length of the n-gram. One of the most commonly used versions is **ROUGE-1**.

**ROUGE-1** measures the overlap of unigrams (single tokens) between the generated and reference texts. As a recall-focused metric (hence the "R" in ROUGE), it assesses how much of the reference text is captured in the generated output.

Recall is the ratio of matching tokens to the total number of tokens in the reference text:

$$\text{recall} \stackrel{\text{def}}{=} \frac{\text{Number of matching tokens}}{\text{Total number of tokens in reference texts}}$$

Formally, ROUGE-1 is defined as:

$$\text{ROUGE-1} \stackrel{\text{def}}{=} \frac{\sum_{(g,r) \in \mathcal{D}} \sum_{t \in r} \text{count}(t, g)}{\sum_{(g,r) \in \mathcal{D}} \text{length}(r)}$$

Here,  $\mathcal{D}$  is the dataset of (generated text, reference text) pairs,  $\text{count}(t, g)$  counts how often a token  $t$  from the reference text  $r$  appears in the generated text  $g$ , and the denominator is the total token count across all reference texts.

To understand this calculation, consider a simple example:

Reference text	Generated text
Large language models are very important for text processing.	Large language models are useful in processing text.

Let's use word-level tokenization and calculate:

- **Matching words:** large, language, models, are, processing, text (6 words)
- **Total words in the reference text:** 9
- **ROUGE-1:**  $\frac{6}{9} \approx 0.67$

A ROUGE-1 score of 0.67 means roughly two-thirds of the words from the reference text appear in the generated text. However, this number alone has little value. ROUGE scores are only useful for *comparing* how different language models perform on the same test set, as they indicate which model more effectively captures the content of the reference texts.

**ROUGE-N** extends the ROUGE metric from unigrams to n-grams while using the same formula.

**ROUGE-L** relies on the **longest common subsequence (LCS)**. This identifies the longest sequence of tokens that appear in both the generated and reference texts in the same order, without needing to be adjacent.

Let  $L_g$  be the length of the generated text, and  $L_r$  the length of the reference text. Then:

$$\text{recall}_{\text{LCS}} \stackrel{\text{def}}{=} \frac{\text{LCS}(L_g, L_r)}{L_r}, \quad \text{precision}_{\text{LCS}} \stackrel{\text{def}}{=} \frac{\text{LCS}(L_g, L_r)}{L_g}$$

Here,  $\text{LCS}(L_g, L_r)$  represents the number of tokens in the LCS between the generated text  $g$  and the reference text  $r$ . The **recall** measures the proportion of the reference text captured by the LCS, while the **precision** measures the proportion of the generated text that matches the reference. Recall and precision are combined into a single metric as follows:

$$\text{ROUGE-L} \stackrel{\text{def}}{=} (1 + \beta^2) \times \frac{\text{recall}_{\text{LCS}} \times \text{precision}_{\text{LCS}}}{\text{recall}_{\text{LCS}} + \beta^2 \times \text{precision}_{\text{LCS}}}$$

Here,  $\beta$  controls the trade-off between precision and recall in the ROUGE-L score. Since ROUGE favors recall,  $\beta$  is usually set high. In Chin-Yew Lin's 2004 paper<sup>6</sup>,  $\beta$  was set to 8.

---

<sup>6</sup> Chin-Yew Lin, "ROUGE: A Package for Automatic Evaluation of Summaries," Proceedings of the Workshop on Text Summarization Branches Out, 2004.

Let's revisit the two texts used to illustrate ROUGE-1. For these sentences, there are two valid longest common subsequences, each with a length of 5 words:

LCS 1	LCS 2
Large, language, models, are, text	Large, language, models, are, processing

Both subsequences are the longest sequences of tokens that appear in the same order in both sentences, but not necessarily consecutively. When multiple LCS options exist, ROUGE-L can use any of them since their lengths are identical.

Here's how the calculations work here. The length of the LCS is 5 words. The reference text is 9 words long, and the generated text is 8 words long. Thus, recall and precision are:

$$\text{recall}_{\text{LCS}} = \frac{5}{9} \approx 0.56, \quad \text{precision}_{\text{LCS}} = \frac{5}{8} \approx 0.63$$

With  $\beta = 8$ , ROUGE-L is then given by:

$$\text{ROUGE-L} = \frac{(1 + 8^2) \cdot 0.56 \cdot 0.63}{0.56 + 8^2 \cdot 0.63} \approx 0.56$$

ROUGE scores range from 0 to 1, where 1 means a perfect match between generated and reference texts. However, even excellent summaries or translations rarely approach 1 in practice.

Choosing the right ROUGE metric depends on the task:

- ROUGE-1 and ROUGE-2 are standard starting points. ROUGE-1 checks overall content similarity using unigram overlap, while ROUGE-2 evaluates local fluency and phrase accuracy using bigram matches.
- For tasks like summarization or paraphrasing, where word order can vary, ROUGE-L is helpful. It detects similarity despite rewording, making it suitable for linguistic variation.
- In cases where preserving longer patterns is key—such as maintaining technical terms or idioms—higher-order metrics like ROUGE-3 or ROUGE-4 might be more relevant.

A combination of metrics, such as ROUGE-1, ROUGE-2, and ROUGE-L, often gives a more balanced evaluation, covering both content overlap and structural flexibility.

Keep in mind, though, that ROUGE has limitations. It measures lexical overlap but not semantic similarity or factual correctness. To address these gaps, ROUGE is often paired with human evaluations or other methods for a fuller assessment of text quality.

### 2.6.3. Human Evaluation

Automated metrics are useful, but human evaluation is still necessary to assess language models. Humans can evaluate qualities that automated metrics often miss, like fluency and accuracy. Two common approaches for human evaluation are Likert scale ratings and Elo ratings.

**Likert scale ratings** involve assigning scores to outputs using a fixed, typically **symmetric scale**. Raters judge the quality by selecting a score, often from  $-2$  to  $2$ , where each scale point corresponds to a descriptive label. For instance,  $-2$  could mean "Strongly Disagree" or "Poor," while  $2$  might mean "Strongly Agree" or "Excellent." The scale is symmetric because it includes equal levels of agreement and disagreement around a neutral midpoint, making positive and negative responses easier to interpret.

Likert scales are flexible for evaluating different aspects of language model outputs, such as fluency, coherence, relevance, and accuracy. For example, a rater could separately rate a sentence's grammatical correctness and its relevance to a prompt, both on a scale from  $-2$  to  $2$ .

However, the method has limitations. One issue is **central tendency bias**, where some raters avoid extreme scores and stick to the middle of the scale. Another challenge is inconsistency in how raters interpret the scale—some may reserve a  $2$  for exceptional outputs, while others may assign it to any high-quality response.

To mitigate these issues, researchers often involve multiple raters, phrase similar questions in different ways for the same rater, and use detailed rubrics that clearly define each scale point.

Let's illustrate Likert scale evaluation using a scenario where machine-generated summaries of news articles are assessed.

Human raters compare a model-generated summary to the original article. They rate it on three aspects using 5-point Likert scales: coherence, informativeness, and factual accuracy.

For example, consider the news article on the left and the generated summary on the right:

## CIRCLES THE WORLD REPORTS 'I FEEL FINE'

MOSCOW—(AP-UPI)—The Russians rocketed the first man into space today and brought him back safely yesterday, the Soviet Union announced. A youthful family man, Maj. Yuri Gagarin, father of two children, was hurtled nearly 200 miles above the earth and sent hurtling around it at the rate of once every 89 minutes. From inside his lonely cabin Gagarin radioed: "I feel fine."

Soviet scientists could see him on their television screens as Gagarin felt himself go weightless at the moment of lift-off. He was alone in the tiny capsule.

When he landed, Tass reported he said: "I feel well. I have no injuries or bruises."

Gagarin was in space for 108 minutes, meaning he completed just slightly more than one orbit of the earth.

"GAGARIN IS THE GREATEST ACHIEVEMENT!"

The fact signalled man's first conquest of space, and a noted British scientist at once hailed it as "the greatest scientific achievement in the history of man." Eventually it may open the planets to exploration by man of all nations.

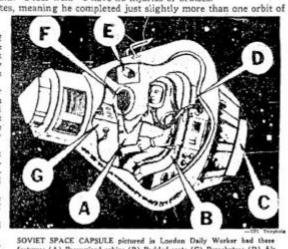
The response from Soviet Premier Khrushchev, although brief, was a message of congratulations to Gagarin he said the "entire Soviet people acclaim your valiant and heroic act." The entire world, he added, centuries as an example of courage, gallantry and heroism in the name of mankind.

Astronauts, sent from space, were read to spell-bound radio listeners glued to their sets. The capsule, which had been bound in a protective cover, was carried in a snow-covered Moscow square.

Three cameras recorded Gagarin's sense of well-being despite the shock of blast-off and his following state of weightlessness. Gagarin, who was described as a "wild duck," told Tass: "I feel fine."

Tass, the official Soviet news agency, announced: "Moscow Mail" Yuri Gagarin safely landed in the prearranged area of the USSR."

The launching and successful return of a human to earth gave Russia victory in the grueling race with the U.S. to put the first



SOVIET SPACE CAPSULE pictured in London Daily Worker had these features: (A) Pressurized cabin; (B) Padded seat; (C) Parachutes; (D) Air supply; (E) TV camera, microphones; (F) Porthole; (G) Instrument panel.

EICHMANN TOLD

This is a historic news article reporting on Yuri Gagarin becoming the first human in space on April 12, 1961. The article details how the Soviet Union successfully launched Gagarin, described as a "youthful family man" and father of two, into orbit around Earth. He spent 108 minutes in space, completing slightly more than one orbit, reaching an altitude of nearly 200 miles.

During the flight, Gagarin radioed "I feel fine" and was observed on television screens by Soviet scientists as he experienced weightlessness. Upon landing safely at a pre-arranged location, he reported having no injuries.

The achievement was hailed as "the greatest scientific achievement in the history of man" by a British scientist, and Soviet Premier Khrushchev praised it as a feat that would be remembered for centuries. The Soviet public followed the event closely, gathering around radios in their homes and in Moscow squares to hear updates.

The article includes a diagram from the London Daily Worker showing various features of the Soviet space capsule, including the pressurized cabin, padded seat, parachutes, air supply, TV cameras, microphone, porthole, and instrument panel.

This accomplishment represented a significant victory for the Soviet Union in its space race with the United States to put the first human in space. The article also notes an interesting detail that Gagarin's name means "wild duck" in English.

Raters assess the summary based on how effectively it meets these three criteria.

The scale for assessing **coherence**—that is, how well-organized, readable, and logically connected the summary is—might look like this:

Very poor	Poor	Acceptable	Good	Excellent
-2	-1	0	1	2

The scale for assessing **informativeness**, that is, how well the summary captures the essence and main points of the original article, might look this way:

Not informative	Slightly	Moderately	Very	Extremely
-2	-1	0	1	2

The scale for assessing **factual accuracy**—that is, how precisely the summary represents facts and data from the original article—might look like this:

Very low	Some inaccuracies	Mostly Accurate	Very accurate	Perfect
-2	-1	0	1	2

In this example, raters would select one option for each of the three aspects. The use of descriptive anchors at each point on the scale helps standardize understanding among raters.

After collecting ratings from multiple raters across various summaries, researchers analyze the data through several approaches:

- Calculate average scores for each aspect across all summaries to get an overall performance measure.
- Compare scores across different versions of the model to track improvements.
- Analyze the correlation between different aspects (e.g., is high coherence associated with high factual accuracy?).

Although Likert scale ratings were originally intended for humans, the rise of advanced **chat LMs** means raters can now be either humans or language models.

**Pairwise comparison** is a method where two outputs are evaluated side-by-side, and the better one is chosen based on specific criteria. This simplifies decision-making, especially when outputs are of similar quality or changes are minor.

The method builds on the principle that relative judgments are easier than absolute ones. Binary choices often produce more consistent and reliable results than absolute ratings.

In practice, raters compare pairs of outputs, such as translations, summaries, or answers, and decide which is better based on criteria like coherence, informativeness, or factual accuracy.

For example, in machine translation evaluation, raters compare pairs of translations for the same source sentence, selecting which one better preserves the original meaning in the target language. By repeating this process across many pairs, evaluators can compare different models or model versions.

Pairwise comparison helps rank models or model versions by having each rater evaluate multiple pairs, with each model compared against others several times. This repetition minimizes individual

biases, resulting in more reliable evaluations. A related approach is **ranking**, where a rater orders multiple responses by quality. Ranking requires less effort than pairwise comparisons while still capturing relative quality.

Results from pairwise comparisons are typically analyzed statistically to determine significant differences between models. A common method for this analysis is the Elo rating system.

**Elo ratings**, originally created by Arpad Elo for ranking chess players, can be adapted for language model evaluation. The system assigns ratings based on “wins” and “losses” in direct comparisons, quantifying relative model performance.

In language model evaluation, all models typically start with an initial rating, often set to 1500. When two models are compared, the probability of one model “winning” is calculated using their current ratings. After each comparison, their ratings are updated based on the actual result versus the expected result.

The probability of model  $A$  with rating  $\text{Elo}(A)$  winning against model  $B$  with rating  $\text{Elo}(B)$  is:

$$\Pr(A \text{ wins}) = \frac{1}{1 + 10^{(\text{Elo}(B) - \text{Elo}(A))/400}}$$

The value 400 in the Elo formula acts as a scaling factor, creating a logarithmic relationship between rating differences and winning probabilities. Arpad Elo chose this number ensuring that a 400-point rating difference reflects 10:1 odds in favor of the higher-rated chess player. While originally designed for chess, this scaling factor has proven effective in other contexts, including language model evaluation.

After a match, ratings are updated using the formula:

$$\text{Elo}(A) \leftarrow \text{Elo}(A) + k \times (\text{score}(A) - \Pr(A \text{ wins})),$$

where  $k$  (typically between 4 and 32) controls the maximum rating change, and  $\text{score}(A)$  reflects the outcome: 1 for a win, 0 for a loss, and 0.5 for a draw.

Consider an example with three models:  $\text{LM}_1$ ,  $\text{LM}_2$ , and  $\text{LM}_3$ . We’ll evaluate them based on their ability to generate coherent text continuations. Assume their initial ratings are:

$$\begin{aligned}\text{Elo}(\text{LM}_1) &= 1500 \\ \text{Elo}(\text{LM}_2) &= 1500 \\ \text{Elo}(\text{LM}_3) &= 1500\end{aligned}$$

We’ll use  $k = 32$  for this example.

Consider this prompt: “*The scientists were shocked when they discovered...*”

Continuation by  $\text{LM}_1$ : “*...a new species of butterfly in the Amazon rainforest. Its wings were unlike anything they had ever seen before.*”

Continuation by LM<sub>2</sub>: "...that the ancient artifact they unearthed was emitting a faint, pulsating light. They couldn't explain its source."

Continuation by LM<sub>3</sub>: "...the results of their experiment contradicted everything they thought they knew about quantum mechanics."

Let's say we conduct pairwise comparisons and get the following results:

1. LM<sub>1</sub> vs LM<sub>2</sub>: LM<sub>1</sub> wins
  - o  $\text{Pr}(\text{LM}_1 \text{ wins}) = 1/(1 + 10^{((1500-1500)/400)}) = 0.5$
  - o New rating for LM<sub>1</sub>  $\leftarrow 1500 + 32(1 - 0.5) = 1516$
  - o New rating for LM<sub>2</sub>  $\leftarrow 1500 + 32(0 - 0.5) = 1484$
2. LM<sub>1</sub> vs LM<sub>3</sub>: LM<sub>3</sub> wins
  - o  $\text{Pr}(\text{LM}_1 \text{ wins}) = 1/(1 + 10^{((1500-1516)/400)}) \approx 0.523$
  - o New rating for LM<sub>1</sub>  $\leftarrow 1516 + 32(0 - 0.523) \approx 1499$
  - o New rating for LM<sub>3</sub>  $\leftarrow 1500 + 32(1 - 0.477) \approx 1517$
3. LM<sub>2</sub> vs LM<sub>3</sub>: LM<sub>3</sub> wins
  - o  $\text{Pr}(\text{LM}_2 \text{ wins}) = 1/(1 + 10^{((1517-1484)/400)}) \approx 0.453$
  - o New rating for LM<sub>2</sub>  $\leftarrow 1484 + 32(0 - 0.453) \approx 1470$
  - o New rating for LM<sub>3</sub>  $\leftarrow 1517 + 32(1 - 0.547) \approx 1531$

Final ratings after these comparisons:

$$\begin{aligned}\text{Elo(LM}_1\text{)} &= 1499 \\ \text{Elo(LM}_2\text{)} &= 1470 \\ \text{Elo(LM}_3\text{)} &= 1531\end{aligned}$$

Elo ratings quantify how models perform relative to each other. In this case, LM<sub>3</sub> is the strongest, followed by LM<sub>1</sub>, with LM<sub>2</sub> ranking last.

Performance isn't judged from a single match. Instead, multiple pairwise matches are used. This limits the effects of random fluctuations or biases in individual comparisons, giving a better estimate of each model's performance.

A variety of prompts or inputs ensures evaluation across different contexts and tasks. When human raters are involved, several raters assess each comparison to reduce individual bias.

To avoid order effects, both the sequence of comparisons and the presentation of outputs are randomized. Elo ratings are updated after every comparison.

How many matches are needed until the results are reliable? There's no universal number that applies to all cases. As a general guideline, some researchers suggest that each model should participate in at least 100–200 comparisons before considering the Elo ratings stable and ideally 500+ comparisons for high confidence. However, for high-stakes evaluations or when comparing very similar models, thousands of comparisons may be necessary.

Statistical methods can be used to calculate a **confidence interval** for a model's Elo rating. Explaining these techniques is beyond the scope of this book. For those interested, the **Bradley-Terry model** and **bootstrap resampling** are good starting points. Both are well-documented, with resources linked on the book's wiki.

Elo ratings provide a continuous scale for ranking models, making it easier to track incremental improvements. The system rewards wins against strong opponents more than wins against weaker ones, and it can handle incomplete comparison data, meaning not every model needs to be compared against every other model. However, the choice of  $k$  significantly affects rating volatility; a poorly chosen  $k$  can undermine the evaluation's stability.

To address these limitations, Elo ratings are often used alongside other evaluation methods. For instance, researchers might use Elo ratings for ranking models in pairwise comparisons, while collecting Likert scale ratings to assess absolute quality. This combined approach yields a more comprehensive view of a language model's performance.

Now that we've covered language modeling and evaluation methods, let's explore a more advanced model architecture: recurrent neural networks (RNNs). RNNs made significant progress in processing text. They introduced the ability to maintain context over long sequences, allowing for the creation of more powerful language models.

## Chapter 3. Recurrent Neural Network

In this chapter, we discuss a key neural network architecture that changed how machines handle sequences: the recurrent neural network. We'll cover its structure and use in language modeling. While recurrent neural networks are largely superseded by transformers and other attention-based architectures in many modern applications, they remain foundational to understanding how machine learning approaches sequential data.

### 3.1. Elman RNN

A **recurrent neural network**, or **RNN**, is a neural network designed for sequential data. Unlike **feedforward neural networks**, RNNs include loops in their connections, enabling information to carry over from one step in the sequence to the next. This makes them well-suited for tasks like time series analysis, natural language processing, and other sequential data problems.

To illustrate the sequential nature of RNNs, let's consider a neural network with a single unit. Consider the input document "*Learning from text is cool.*" Ignoring case and punctuation, the matrix representing this document would be as follows:

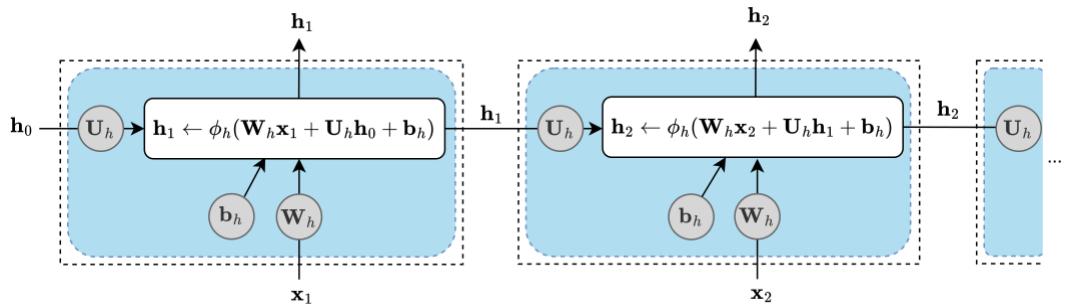
Word	Embedding vector
learning	[0.1,0.2,0.6] <sup>T</sup>
from	[0.2,0.1,0.4] <sup>T</sup>
text	[0.1,0.3,0.3] <sup>T</sup>
is	[0.0,0.7,0.1] <sup>T</sup>
cool	[0.5,0.2,0.7] <sup>T</sup>
PAD	[0.0,0.0,0.0] <sup>T</sup>

Each row of the matrix represents a word's embedding learned during neural network training. The order of words is preserved. The matrix dimensions are (sequence length, embedding dimensionality). Sequence length specifies the maximum number of words in a document. Shorter documents are padded with padding tokens, while longer ones are truncated. **Padding** uses dummy embeddings, usually **zero vectors**.

More formally, the matrix would look like this:

$$\mathbf{X} = \begin{bmatrix} 0.1 & 0.2 & 0.6 \\ 0.2 & 0.1 & 0.4 \\ 0.1 & 0.3 & 0.3 \\ 0.0 & 0.7 & 0.1 \\ 0.5 & 0.2 & 0.7 \\ 0.0 & 0.0 & 0.0 \end{bmatrix}$$

Here, we have five 3D embedding vectors,  $\mathbf{x}_1, \dots, \mathbf{x}_5$ , representing each word in the document. For instance,  $\mathbf{x}_1 = [0.1,0.2,0.6]^T$ ,  $\mathbf{x}_2 = [0.2,0.1,0.4]^T$ , and so on. The sixth vector is a padding vector. The single-unit RNN used to process this sequence is structured as follows:



The same unit receives as input a sequence of embedding vectors, one at a time, and outputs, at each time step  $t$ , the hidden state  $\mathbf{h}_t$ . This unit is known as the **Elman RNN**, named after Jeffrey Locke Elman, who introduced the **simple recurrent neural network** in 1990.

Unlike a unit in a multilayer perceptron, which takes a vector and returns a scalar as shown in Figure 1.1, an RNN unit returns a vector, functioning as an entire layer.

As shown, an RNN receives two inputs at each time step  $t$ : an input embedding vector  $\mathbf{x}_t$  (typically a word embedding) and a hidden state vector  $\mathbf{h}_{t-1}$  from the previous time step. It outputs an updated hidden state  $\mathbf{h}_t$ . This characteristic—using its own output from the previous time step as an input—gives the network its “recurrent” nature. The initial hidden state  $\mathbf{h}_0$  is usually initialized with a **zero vector**.

A **hidden state** is a vector that stores information from all previous time steps in a sequence. It acts as the network’s “memory,” enabling past information to influence future predictions. At each time step, the hidden state is updated using the current input and the previous hidden state. This is critical for sequential tasks like language modeling, where context from earlier words helps predict the next word.

To deepen the network, we add a second RNN layer. The first layer’s outputs,  $\mathbf{h}_t$ , become inputs to the second, whose outputs are the network’s final outputs:

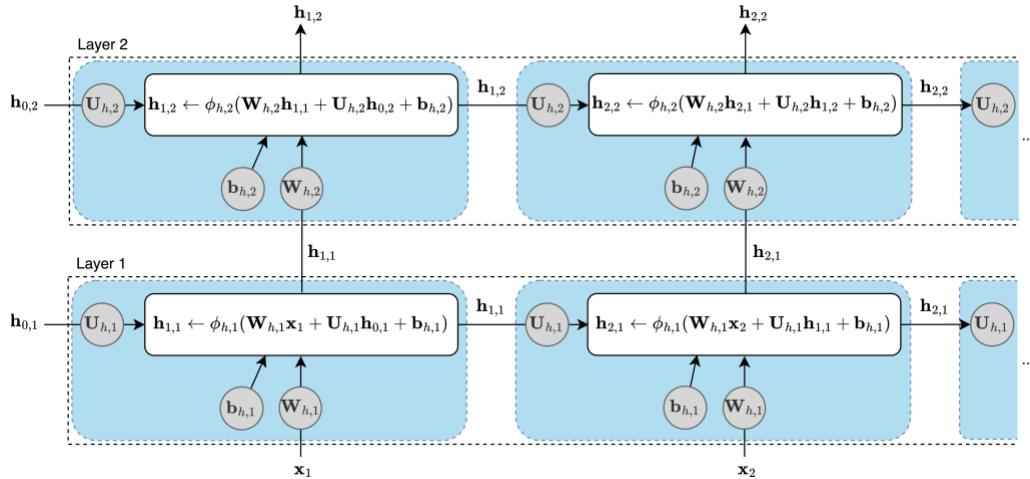


Figure 3.1: A two-layer Elman RNN. The first layer’s outputs serve as inputs to the second layer.

## 3.2. Mini-Batch Gradient Descent

Before coding the RNN model, we need to discuss the shape of the input data. In Section 1.7, we used the entire dataset for each gradient descent step. Here, and for training all future models, we’ll adopt **mini-batch gradient descent**, a widely used method for large models and datasets. Mini-batch gradient descent calculates derivatives over smaller data subsets, which speeds up learning and reduces memory usage.

With mini-batch gradient descent, the data shape is organized as (batch size, sequence length, embedding dimensionality). This structure divides the training set into fixed-size mini-batches, each containing sequences of embeddings with consistent lengths. (From this point on, “batch” and “mini-batch” will be used interchangeably.)

For example, if the batch size is 2, the sequence length is 4, and the embedding dimensionality is 3, the mini-batch can be represented as:

$$\text{batch}_1 = \begin{bmatrix} \text{seq}_{1,1} & \text{seq}_{1,2} & \text{seq}_{1,3} & \text{seq}_{1,4} \\ \text{seq}_{2,1} & \text{seq}_{2,2} & \text{seq}_{2,3} & \text{seq}_{2,4} \end{bmatrix}$$

Here,  $\text{seq}_{i,j}$ , for  $i \in \{1,2\}$  and  $j \in \{1, \dots, 4\}$  is an embedding vector.

For example, if we have the following embeddings for each sequence:

$$\text{seq}_1: \begin{bmatrix} [0.1, 0.2, 0.3] \\ [0.4, 0.5, 0.6] \\ [0.7, 0.8, 0.9] \\ [1.0, 1.1, 1.2] \end{bmatrix}$$

$$\text{seq}_2: \begin{bmatrix} [1.3, 1.4, 1.5] \\ [1.6, 1.7, 1.8] \\ [1.9, 2.0, 2.1] \\ [2.2, 2.3, 2.4] \end{bmatrix}$$

The mini-batch will look like this:

$$\text{batch}_1 = \begin{bmatrix} [0.1, 0.2, 0.3] & [0.4, 0.5, 0.6] & [0.7, 0.8, 0.9] & [1.0, 1.1, 1.2] \\ [1.3, 1.4, 1.5] & [1.6, 1.7, 1.8] & [1.9, 2.0, 2.1] & [2.2, 2.3, 2.4] \end{bmatrix}$$

During each step of gradient descent, we:

1. Select a mini-batch from the training set,
2. Pass it through the neural network,
3. Compute the loss,
4. Calculate gradients,
5. Update model parameters,
6. Repeat from step 1.

Mini-batch gradient descent often achieves faster **convergence** compared to using the entire training set per step. It efficiently handles large models and datasets by using modern hardware's parallel processing capabilities. In PyTorch, models require the first dimension of the input data to be the batch dimension, even if there's only one example in the batch.

### 3.3. Programming an RNN

Let's implement an Elman RNN unit:

```
import torch
import torch.nn as nn

class ElmanRNNUnit(nn.Module):
    def __init__(self, emb_dim):
        super().__init__()
        self.Uh = nn.Parameter(torch.randn(emb_dim, emb_dim)) ①
        self.Wh = nn.Parameter(torch.randn(emb_dim, emb_dim)) ②
        self.b = nn.Parameter(torch.zeros(emb_dim)) ③

    def forward(self, x, h):
        return torch.tanh(x @ self.Wh + h @ self.Uh + self.b) ④
```

In the constructor:

- Lines ① and ② initialize `self.Uh` and `self.Wh`, the weight matrices for the hidden state and input vector, with random values.
- Line ③ sets `self.b`, the bias vector, to zero.

In the `forward` method, line ❸ handles the computation for each time step. It processes the current input  $x$  and the previous hidden state  $h$ , both shaped  $(batch\_size, emb\_dim)$ , combines them with the weight matrices and bias, and applies the `tanh` activation. The output is the new hidden state, also of shape  $(batch\_size, emb\_dim)$ .

The @ the **matrix multiplication** operator in PyTorch. We use  $x @ self.Wh$  rather than `self.Wh @ x` because of the way PyTorch handles batch dimensions in matrix multiplication. When working with batched inputs,  $x$  has a shape of  $(batch\_size, emb\_dim)$ , while `self.Wh` has a shape of  $(emb\_dim, emb\_dim)$ . Remember from Section 1.6 that for two matrices to be multipliable, the number of columns in the left matrix must be the same as the number of rows in the right matrix. This is satisfied in  $x @ self.Wh$ .

Now, let's define the class `ElmanRNN`, which implements a two-layer Elman RNN using `ElmanRNNUnit` as its core building block:

```
class ElmanRNN(nn.Module):
    def __init__(self, emb_dim, num_layers):
        super().__init__()
        self.emb_dim = emb_dim
        self.num_layers = num_layers
        self.rnn_units = nn.ModuleList(
            [ElmanRNNUnit(emb_dim) for _ in range(num_layers)])
    ) ❶

    def forward(self, x):
        batch_size, seq_len, emb_dim = x.shape ❷
        h_prev = [
            torch.zeros(batch_size, emb_dim, device=x.device) ❸
            for _ in range(self.num_layers)
        ]
        outputs = []
        for t in range(seq_len): ❹
            input_t = x[:, t]
            for l, rnn_unit in enumerate(self.rnn_units):
                h_new = rnn_unit(input_t, h_prev[l])
                h_prev[l] = h_new  # Update hidden state
                input_t = h_new  # Input for next layer
            outputs.append(input_t) ❽ Collect outputs
        return torch.stack(outputs, dim=1) ❾
```

In line ❶ of the constructor, we initialize the RNN layers by creating a `ModuleList` containing `ElmanRNNUnit` instances—one per layer. Using `ModuleList` instead of a regular Python list ensures the parent module (`ElmanRNN`) properly registers all RNN unit parameters. This guarantees that calling `.parameters()` or `.to(device)` on the parent module includes parameters from all modules in the `ModuleList`.

In the `forward` method:

- Line ❷ extracts `batch_size`, `seq_len`, and `emb_dim` from the input tensor `x`.
- Line ❸ initializes the hidden states `h_prev` for all layers with zero tensors. Each hidden state in the list has the shape `(batch_size, emb_dim)`.

We store hidden states for each layer in a list instead of a multidimensional tensor because we need to modify them during processing. In-place modifications of tensors can disrupt PyTorch's automatic differentiation system, which might result in incorrect gradient calculations.

- Line ❹ iterates over time steps `t` in the input sequence. For each `t`:
  - Extract the input at time `t`: `input_t = x[:, t]`.
  - For each layer `l`:
    - Compute the new hidden state `h_new` from `input_t` and `h_prev[1]`.
    - Update the hidden state: `h_prev[1] = h_new` (updates in place).
    - Set `input_t = h_new` to pass to the next layer.
  - Append the output of the last layer: `outputs.append(input_t)`.
- Once all time steps are processed, line ❺ converts the `outputs` list into a tensor by stacking it along the time dimension. The resulting tensor has the shape `(batch_size, seq_len, emb_dim)`.

### 3.4. RNN as a Language Model

An RNN-based language model uses ElmanRNN as its building block:

```
class RecurrentLanguageModel(nn.Module):
    def __init__(self, vocab_size, emb_dim, num_layers, pad_idx):
        super().__init__()
        self.embedding = nn.Embedding(
            vocab_size,
            emb_dim,
            padding_idx=pad_idx
        ) ❶
        self.rnn = ElmanRNN(emb_dim, num_layers)
        self.fc = nn.Linear(emb_dim, vocab_size)

    def forward(self, x):
        embeddings = self.embedding(x)
        rnn_output = self.rnn(embeddings)
        logits = self.fc(rnn_output)
        return logits
```

The `RecurrentLanguageModel` class integrates three components: an embedding layer, the `ElmanRNN` defined earlier, and a final linear layer.

In the constructor, line ❶ defines the embedding layer. This layer transforms input token indices into dense vectors. The `padding_idx` parameter ensures that padding tokens are represented by zero vectors. (We'll cover the embedding layer in the next section.)

Next, we initialize the custom `ElmanRNN`, specifying the embedding dimensionality and the number of layers. Finally, we add a fully connected layer, which converts the RNN's output into vocabulary-sized logits for each token in the sequence.

In the `forward` method:

- We pass the input `x` through the embedding layer. Input `x` has shape `(batch_size, seq_len)`, and the output `embeddings` have shape `(batch_size, seq_len, emb_dim)`.
- We then pass the embedded input through our `ElmanRNN`, obtaining `rnn_output` with shape `(batch_size, seq_len, emb_dim)`.
- Finally, we apply the fully connected layer to the RNN output, producing logits for each token in the vocabulary at each position in the sequence. The output logits have shape `(batch_size, seq_len, vocab_size)`.

## 3.5. Embedding Layer

An **embedding layer**, implemented as `nn.Embedding` in PyTorch, maps token indices from a vocabulary to dense, fixed-size vectors. It acts as a learnable lookup table, where each token is assigned a unique embedding vector. During training, these vectors are adjusted to capture meaningful numerical representations of the tokens.

Here's an example to show how an embedding layer works. Imagine a vocabulary with five tokens, indexed from 0 to 4. We want each token to have a three-dimensional embedding vector. To begin, we create an embedding layer:

```
import torch
import torch.nn as nn

vocab_size = 5 # Number of unique tokens
emb_dim = 3    # Size of each embedding vector
emb_layer = nn.Embedding(vocab_size, emb_dim)
```

The embedding layer initializes the embedding matrix **E** with random values. In this case, the matrix has 5 rows (one for each token) and 3 columns (the embedding dimensionality):

$$\mathbf{E} = \begin{bmatrix} 0.2 & -0.4 & 0.1 \\ -0.3 & 0.8 & -0.5 \\ 0.7 & 0.1 & -0.2 \\ -0.6 & 0.5 & 0.4 \\ 0.9 & -0.7 & 0.3 \end{bmatrix}$$

Each row in **E** represents the embedding for a specific token in the vocabulary.

Now, let's input a sequence of token indices:

```
token_indices = torch.tensor([0, 2, 4])
```

The embedding layer retrieves the rows of  $\mathbf{E}$  corresponding to the input indices:

$$\text{Embeddings} = \begin{bmatrix} 0.2 & -0.4 & 0.1 \\ 0.7 & 0.1 & -0.2 \\ 0.9 & -0.7 & 0.3 \end{bmatrix}$$

This output is a matrix whose number of rows equals the input sequence length and whose number of columns equals the embedding dimensionality:

```
embeddings = embedding_layer(token_indices)
print(embeddings)
```

The output might look like this:

```
tensor([[ 0.2, -0.4,  0.1],
       [ 0.7,  0.1, -0.2],
       [ 0.9, -0.7,  0.3]])
```

The embedding layer can manage padding tokens as well. Padding ensures sequences in a mini-batch have the same length. To prevent the model from updating embeddings for padding tokens during training, the layer maps them to a zero vector that remains unchanged. For example, if we define the padding index:

```
emb_layer = nn.Embedding(vocab_size, emb_dim, padding_idx=0)
```

The embedding for token 0 (padding token) is always  $[0,0,0]^T$ .

Given the input:

```
token_indices = torch.tensor([0, 2, 4])
embeddings = emb_layer(token_indices)
print(embeddings)
```

The result would be:

```
tensor([[ 0.0,  0.0,  0.0], # Padding token
       [ 0.7,  0.1, -0.2], # Token 2 embedding
       [ 0.9, -0.7,  0.3]]) # Token 4 embedding
```

With modern language models, vocabularies often include hundreds of thousands of tokens, and embedding dimensions are typically around a thousand. This makes the embedding matrix a significant part of the model, sometimes containing nearly 2 billion parameters.

## 3.6. Training a Language Model

Start by importing libraries and defining utility functions:

```

import torch, torch.nn as nn

def set_seed(seed):
    random.seed(seed)
    torch.manual_seed(seed)
    torch.cuda.manual_seed_all(seed) ①
    torch.backends.cudnn.deterministic = True ②
    torch.backends.cudnn.benchmark = False ③

```

The `set_seed` function enforces **reproducibility** by setting random seeds. It sets the Python random seed, the PyTorch CPU seed, and, in line ①, the CUDA seed for all GPUs (Graphics Processing Units). CUDA is NVIDIA’s parallel computing platform and API that enables significant performance improvements in computing by leveraging the power of GPUs. Using `torch.cuda.manual_seed_all` ensures consistent GPU-based random behavior, while lines ② and ③ disable CUDA’s auto-tuner and enforce deterministic algorithms, guaranteeing identical results across different GPU models.

With the model class ready, we’ll train our neural language model. First, we install the `transformers` package—an open-source library providing APIs and tools to easily download, train and use pretrained models from the **Hugging Face Hub**:

```
$ pip3 install transformers
```

The package offers a Python API for training that works with both **PyTorch** and **TensorFlow**. For now, we only need it to get a tokenizer.

Now we import `transformers`, set the tokenizer, define the hyperparameter values, prepare the data, and instantiate the model, loss function, and optimizer objects:

```

from transformers import AutoTokenizer

device = torch.device("cuda" if torch.cuda.is_available() else "cpu") ①
tokenizer = AutoTokenizer.from_pretrained(
    "microsoft/Phi-3.5-mini-instruct"
) ②
vocab_size = len(tokenizer) ③

emb_dim, num_layers, batch_size, learning_rate, num_epochs = get_hyperparameters()

data_url = "https://www.thelmbbook.com/data/news"
train_loader, test_loader = download_and_prepare_data(
    data_url, batch_size, tokenizer) ④

model = RecurrentLanguageModel(
    vocab_size, emb_dim, num_layers, tokenizer.pad_token_id
)

```

```

initialize_weights(model) ❶
model.to(device)

criterion = nn.CrossEntropyLoss(ignore_index=tokenizer.pad_token_id) ❷
optimizer = torch.optim.AdamW(model.parameters(), lr=learning_rate)

```

Line ❶ detects a CUDA device if it's available. Otherwise, it defaults to CPU.

Most models on the Hugging Face Hub include the tokenizer that was used to train them. Line ❷ initializes the **Phi 3.5 mini** tokenizer. It was trained on a large text corpus using the **byte-pair encoding** algorithm and has a vocabulary size of 32,064.

Line ❸ retrieves the tokenizer's vocabulary size. Line ❹ downloads and prepares the dataset—a collection of news sentences from online articles—tokenizing them and creating `DataLoader` objects that iterate over batches.

Line ❺ initializes the model parameters. Initial parameter values can greatly influence the training process. They can affect how quickly training progresses and the final loss value. Certain initialization techniques, like **Xavier initialization**, have shown good results in practice. The `initialize_weights` function, implementing this method, is defined in the notebook.

Line ❻ creates the loss function with the `ignore_index` parameter. This ensures the loss is not calculated for padding tokens.

Now, let's look at the training loop:

```

for epoch in range(num_epochs): ❶
    model.train() ❷
    for batch in train_loader: ❸
        input_seq, target_seq = batch
        input_seq = input_seq.to(device) ❹
        target_seq = target_seq.to(device) ❺
        batch_size_current, seq_len = input_seq.shape ❻
        optimizer.zero_grad()
        output = model(input_seq)
        output = output.reshape(batch_size_current * seq_len, vocab_size) ❼
        target = target_seq.reshape(batch_size_current * seq_len) ➋
        loss = criterion(output, target) ❾
        loss.backward()
        optimizer.step()

```

Line ❶ iterates over epochs. An **epoch** is a single pass through the entire dataset. Training for multiple epochs can improve the model, especially with limited training data. The number of epochs is a **hyperparameter** that you adjust based on the model's performance on the test set.

Line ❷ calls `model.train()` at the start of each epoch to set the model in training mode. This is important for models that have layers behaving differently during training vs. **evaluation**.

Although our RNN model doesn't use such layers, calling `model.train()` ensures the model is properly configured for training. This avoids unexpected behavior and keeps consistency, especially if future changes add layers dependent on the mode.

Line ❸ iterates over batches. Each batch is a tuple: one tensor contains input sequences, and the other contains target sequences. Lines ❹ and ❺ move these tensors to the same device as the model. If the model and data are on different devices, PyTorch raises an error.

Line ❻ retrieves the batch size and sequence length from `input_seq` (`target_seq` has the same shape). These dimensions are needed to reshape the model's output tensor (`batch_size_current`, `seq_len`, `vocab_size`) and target tensor (`batch_size_current`, `seq_len`) into compatible shapes for the **cross-entropy** loss function. In line ❼, the output is reshaped to `(batch_size_current * seq_len, vocab_size)`, and in line ❽, the target is flattened to `batch_size_current * seq_len`, allowing the loss calculation in line ❾ to process all tokens in the batch simultaneously and return the average loss per token.

## 3.7. Training Data and Loss Computation

When studying neural language models, a key aspect is understanding the structure of a training example. The text corpus is split into overlapping input and target sequences. Each input sequence aligns with a target sequence shifted by one token. This setup trains the model to predict the next word at each position in the sequence.

For instance, take the sentence "*We train a recurrent neural network as a language model.*" After tokenizing it with the Phi 3.5 mini tokenizer, we get:

```
["_We", "_train", "_a", "_rec", "urrent", "_neural", "_network", "_as", "_a",
"_language", "_model", "."]
```

To create one training example, we convert the sentence into input and target sequences by shifting tokens forward by one position:

```
Input: ["_We", "_train", "_a", "_rec", "urrent", "_neural", "_network", "_a
s", "_a", "_language", "_model"]
Target: ["_train", "_a", "_rec", "urrent", "_neural", "_network", "_as", "_a
", "_language", "_model", "."]
```

A training example doesn't need to be a complete sentence. Modern language models process sequences up to their **context window** length—a fixed maximum number of tokens they can handle at once (like 2048, 4096, or 8192 tokens). This context window determines how much text the model can “see” and reason about at any time, which affects its ability to understand relationships between distant parts of text. The training corpus is therefore segmented into chunks matching this context window length, with the target sequence for each chunk shifted forward one position relative to the input.

During training, the RNN processes one token at a time, updating its hidden states layer by layer. At each step, it generates logits aimed at predicting the next token in the sequence. Each logit corresponds to a vocabulary token and is converted into probabilities using **softmax**. These probabilities are then used to compute the loss.

Each sentence results in multiple predictions and losses. For example, the model first processes “\_We” and tries to predict “\_train” by assigning probabilities to all vocabulary tokens. The loss is computed using the probability of “\_train,” as defined in Equation 2.1. Next, the model processes “\_train” to predict “\_a,” generating another loss. This continues for every token in the sequence. In total, the model makes 11 predictions and calculates 11 losses for this example.

The losses are averaged across the tokens in a training example and all examples in the batch. The average loss is then used in backpropagation to update the model’s parameters.

Predicting the next token at each position gives the model many “signals” to learn from, speeding up learning compared to predicting just one hidden token for the whole sequence, as is the case with **masked language models**.

Let’s break down the loss calculation for each position with some made-up numbers:

- **Position 1:**
  - Target token: “\_train”
  - Logit for “\_train”: -0.5
  - After applying softmax to the logits, suppose the probability of “\_train” is 0.1
  - Contribution to the total loss by Equation 2.1 is  $-\log(0.1) = 2.30$
- **Position 2:**
  - Target token: “\_a”
  - Logit for “\_a”: 3.2
  - After softmax, the probability for “\_a”: 0.05
  - Contribution to loss:  $-\log(0.05) = 2.99$
- **Position 3:**
  - The probability for “\_rec”: 0.02
  - Contribution to loss:  $-\log(0.02) = 3.91$
- **Position 4:**
  - The probability for “urrent”: 0.34
  - Contribution to loss:  $-\log(0.34) = 1.08$

We continue until calculating the loss contribution for the final token, the period:

- **Position 11:**
  - Target token: “.”
  - Logit for “.”: -1.2
  - After softmax, the probability for “.”: 0.11
  - Contribution to loss:  $-\log(0.11) = 2.21$

The final loss is calculated by taking the average of these values:

$$\frac{(2.30 + 2.99 + 3.91 + 1.08 + \dots + 2.21)}{11} = 2.11 \text{ (hypothetically)}$$

During training, the objective is to minimize this loss. This involves improving the model so that it assigns higher probabilities to the correct target tokens at each position.

The full code for training the RNN-based language model can be found in [thelmbbook.com/nb/3.1](http://thelmbbook.com/nb/3.1). I used the following hyperparameter values: `emb_dim = 128`, `num_layers = 2`, `batch_size = 128`, `learning_rate = 0.001`, and `num_epochs = 1`.

Here are three continuations for the prompt “*The President*” generated at later training steps:

```
The President refused to comment on the best news in the five on BBC .
The President has been a `` very serious '' and `` unacceptable '' .
The President 's office is not the first time to be able to take the lead .
```

At the start of training, the model generated almost random token sequences. Over time, its outputs improved: it now correctly closes quotes and parentheses in appropriate parts of sentences. Still, the generated continuations remain below the level of advanced LLMs. For instance, the model’s perplexity is 72.41, much higher than the 20 perplexity of the older, relatively small GPT-2 model and far above the perplexity of around 5 achieved by leading LLMs.

This gap has several causes. First, our model is smaller than LLMs, with just 8,292,619 parameters, most of which are in the embedding layer. Second, simple RNN architectures, like the Elman RNN, have clear limitations. While they handle sequential data, they often fail to retain information from earlier tokens as sequences grow. The hidden state gradually “forgets” past inputs. Lastly, RNNs process tokens sequentially, which complicates training of larger models. Each token depends on the processing of the previous one, forcing the GPU to process tokens one at a time rather than leveraging parallel computation.

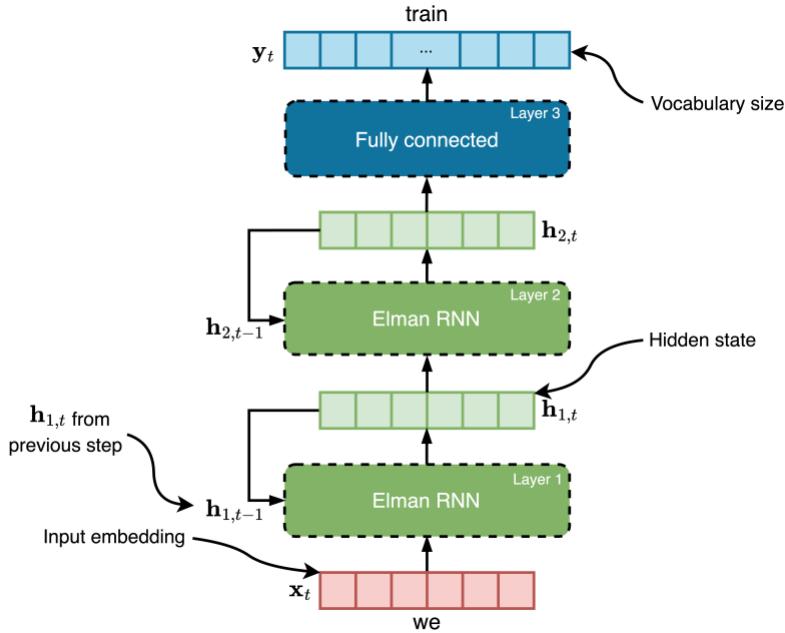
These limitations inspired the development of advanced recurrent architectures like **long short-term memory (LSTM)** networks. LSTMs mitigate some RNN weaknesses but still struggle with very long sequences, such as those spanning thousands of tokens, which are common in modern language models.

The introduction of transformers, discussed in the next chapter, resolved many of these issues. By 2023, transformers have largely replaced RNNs in natural language processing because they handle long-range dependencies better and allow parallel computation.

Interest in RNNs was reignited in 2024 with the invention of the **minLSTM** and **xLSTM** architectures, which achieve performance comparable to Transformer-based models. This resurgence reflects a broader trend in AI research: no model type is ever permanently obsolete. Researchers often revisit and refine older ideas, adapting them to address modern challenges and leverage current hardware capabilities.

### 3.8. Simplified Model Representation

Now that we've covered the math behind language model layers and the structure of the training data, we can simplify the model's representation by representing each unit as a square, just like in Section 1.5. Below is a simplified diagram of the two-layer Elman RNN from Figure 3.1:



Here, we've adjusted the information flow in the diagram from left-to-right, as used in earlier chapters, to bottom-to-top. This is the standard orientation for high-level language model diagrams in the literature. We'll keep this orientation when discussing the Transformer.

With that, we've finished covering recurrent neural networks and the language models built on them. Next, we'll explore transformer neural networks: how they differ from the models we've studied and how they handle tasks like language modeling and document classification.

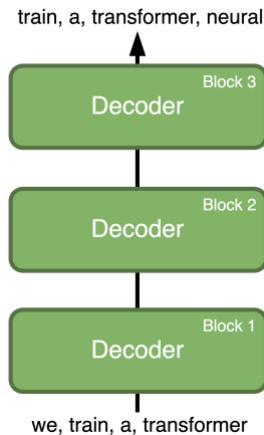
## Chapter 4. Transformer

Transformer models have greatly advanced NLP. They overcome RNNs' limitations in managing long-range dependencies and enable parallel processing of input sequences. There are three main Transformer architectures: encoder-decoder, initially formulated for machine translation; encoder-only, typically used for classification; and decoder-only, commonly found in chat LMs.

In this chapter, we'll explore the decoder-only Transformer architecture in detail, as it is the most widely used approach for training **autoregressive language models**.

The transformer architecture introduces two key innovations: self-attention and positional encoding. Self-attention enables the model to assess how each word relates to all others during prediction, while positional encoding captures word order and sequential patterns. Unlike RNNs, transformers process all tokens simultaneously, using positional encoding to maintain sequential context despite parallel processing. This chapter explores these fundamental elements in detail.

A decoder-only Transformer (referred to simply as “decoder” from here on) is made up of multiple identical<sup>7</sup> layers, known as decoder blocks, stacked vertically as shown below:



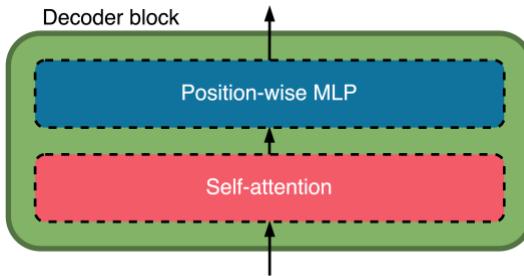
As you can see, training a decoder involves pairing each input sequence with a target sequence that is identical to the input but shifted forward by one token. This approach mirrors the training method used for RNN-based language models.

### 4.1. Decoder Block

Each decoder block has two sub-layers: self-attention and a position-wise multilayer perceptron (MLP) as shown here:

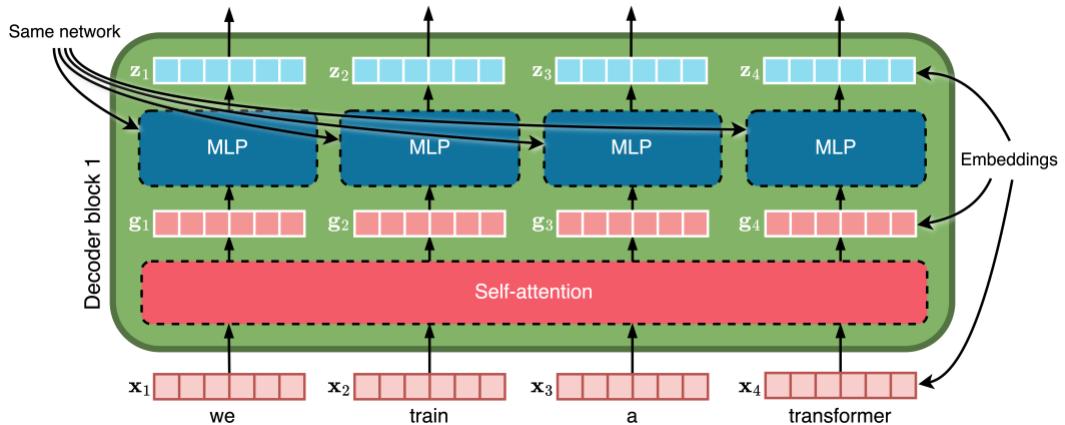
---

<sup>7</sup> Decoder blocks share the same architecture but have distinct trainable parameters unique to each block.



The illustration simplifies certain aspects to avoid introducing too many new concepts at once. We'll introduce the missing details step by step.

Let's take a closer look at what happens in a decoder block, starting with the first one:



The first decoder block processes input token embeddings. For this example, we use 6-dimensional input and output embeddings, though in practice these dimensions grow larger with parameter count and token vocabulary. The **self-attention layer**, transforms each input embedding vector  $\mathbf{x}_t$  into a new vector  $\mathbf{g}_t$  for every token  $t$ , from 1 to  $L$ , where  $L$  represents the input length.

After self-attention, the position-wise MLP independently processes each vector  $\mathbf{g}_t$  one at a time. Each decoder block has its own MLP with unique parameters, and within a block, this same MLP is applied independently to each position's vector, taking one  $\mathbf{g}_t$  as input and producing one  $\mathbf{z}_t$  as output. When the MLP finishes processing each position sequentially, the number of output vectors  $\mathbf{z}_t$  equals the number of input tokens  $\mathbf{x}_t$ .

The output vectors  $\mathbf{z}_t$  then serve as inputs to the next decoder block. This process repeats through each decoder block, preserving the same number of vectors as the input tokens  $\mathbf{x}_t$ .

## 4.2. Self-Attention

To see how **self-attention** works, let's start with an intuitive comparison. Transforming  $\mathbf{g}_t$  into  $\mathbf{z}_t$  is straightforward: a position-wise MLP takes an input vector and outputs a new vector by applying a learned transformation. This is what feedforward networks are designed to do. However, self-attention can seem more complex.

Consider a 5-token example: ["we," "train," "a," "transformer," "model"], and assume a decoder with a maximum input sequence length of 4.

In each decoder block, the self-attention function relies on three tensors of trainable parameters:  $\mathbf{W}^Q$ ,  $\mathbf{W}^K$ , and  $\mathbf{W}^V$ . Here,  $Q$  stands for "query,"  $K$  for "key," and  $V$  for "value."

Let's assume these tensors are  $6 \times 6$ . This means each of the four 6-dimensional input vectors will be transformed into four 6-dimensional output vectors. Let's use the second token,  $\mathbf{x}_2$ , representing the word "train," as our illustrative example. To compute the output  $\mathbf{g}_2$  for  $\mathbf{x}_2$ , the self-attention layer works in six steps.

### 4.2.1. Step 1 of Self-Attention

Compute matrices  $\mathbf{Q}$ ,  $\mathbf{K}$ , and  $\mathbf{V}$  as shown below:

The figure illustrates the computation of query ( $\mathbf{Q}$ ), key ( $\mathbf{K}$ ), and value ( $\mathbf{V}$ ) matrices from the input matrix  $\mathbf{X}$  using weight matrices  $\mathbf{W}^Q$ ,  $\mathbf{W}^K$ , and  $\mathbf{W}^V$ .

- Step 1 (Query):**  $\mathbf{X} \times \mathbf{W}^Q = \mathbf{Q}$
- Step 2 (Key):**  $\mathbf{X} \times \mathbf{W}^K = \mathbf{K}$
- Step 3 (Value):**  $\mathbf{X} \times \mathbf{W}^V = \mathbf{V}$

Each matrix is a 4x6 grid. The input  $\mathbf{X}$  has columns labeled  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4$ . The weight matrices  $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V$  have columns labeled  $\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3, \mathbf{q}_4$ ;  $\mathbf{k}_1, \mathbf{k}_2, \mathbf{k}_3, \mathbf{k}_4$ ; and  $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4$  respectively. The resulting matrices  $\mathbf{Q}, \mathbf{K}, \mathbf{V}$  have columns labeled  $\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3, \mathbf{q}_4$ ;  $\mathbf{k}_1, \mathbf{k}_2, \mathbf{k}_3, \mathbf{k}_4$ ; and  $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4$  respectively.

Figure 4.1: Matrix multiplication in the self-attention layer.

In the illustration, we combined the four input embeddings  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$ , and  $\mathbf{x}_4$  into a matrix  $\mathbf{X}$ . Then, we multiplied  $\mathbf{X}$  by the weight matrices  $\mathbf{W}^Q$ ,  $\mathbf{W}^K$ , and  $\mathbf{W}^V$  to create matrices  $\mathbf{Q}$ ,  $\mathbf{K}$ , and  $\mathbf{V}$ . These matrices hold 6-dimensional query, key, and value vectors, respectively. Since the process

generates the same number of query, key, and value vectors as input embeddings, each input embedding  $\mathbf{x}_t$  corresponds to a query vector  $\mathbf{q}_t$ , a key vector  $\mathbf{k}_t$ , and a value vector  $\mathbf{v}_t$ .

#### 4.2.2. Step 2 of Self-Attention

Taking the second token  $\mathbf{x}_2$  as our example, we compute **attention scores** by taking the dot product of its query vector  $\mathbf{q}_2$  with each key vector  $\mathbf{k}_t$ . Let's assume the resulting scores are:

$$\mathbf{q}_2 \cdot \mathbf{k}_1 = 4.90, \quad \mathbf{q}_2 \cdot \mathbf{k}_2 = 17.15, \quad \mathbf{q}_2 \cdot \mathbf{k}_3 = 9.80, \quad \mathbf{q}_2 \cdot \mathbf{k}_4 = 12.25$$

In vector format:

$$\mathbf{scores}_2 = [4.90, 17.15, 9.80, 12.25]^\top$$

#### 4.2.3. Step 3 of Self-Attention

We now divide each score by the square root of the key vector's dimensionality to obtain the **scaled scores**. In our example, the key vector has a dimensionality of 6, so we divide each score by  $\sqrt{6} \approx 2.45$ , yielding:

$$\mathbf{scaled\_scores}_2 = \left[ \frac{4.9}{2.45}, \frac{17.15}{2.45}, \frac{9.8}{2.45}, \frac{12.25}{2.45} \right]^\top = [2, 7, 4, 5]^\top$$

#### 4.2.4. Step 4 of Self-Attention

We apply the **causal mask** to the scaled scores. (If the reason for using the causal mask isn't clear yet, it will be explained in detail soon.) For the second input position, the causal mask is:

$$\mathbf{causal\_mask}_2 \stackrel{\text{def}}{=} [0, 0, -\infty, -\infty]^\top$$

We add the scaled scores to the causal mask, resulting in the **masked scores**:

$$\mathbf{masked\_scores}_2 = \mathbf{scaled\_scores}_2 + \mathbf{causal\_mask}_2 = [2, 7, -\infty, -\infty]^\top$$

#### 4.2.5. Step 5 of Self-Attention

We apply the **softmax** function to the masked scores to produce the **attention weights**:

$$\mathbf{attention\_weights}_2 = \text{softmax}([2, 7, -\infty, -\infty]^\top)$$

Since scores of  $-\infty$  become zero after applying the exponential function, the attention weights for the third and fourth positions will be zero. The remaining two weights are calculated as:

$$\mathbf{attention\_weights}_2 = \left[ \frac{e^2}{e^2 + e^7}, \frac{e^7}{e^2 + e^7}, 0, 0 \right]^\top \approx [0.0067, 0.9933, 0, 0]^\top$$

Dividing attention scores by the square root of the key dimensionality helps prevent the dot products from growing too large in magnitude as the dimensionality increases,

which could lead to extremely small gradients after applying **softmax** (due to very large negative or positive values pushing the softmax outputs to 0 or 1).

#### 4.2.6. Step 6 of Self-Attention

We compute the output vector  $\mathbf{g}_2$  for the input embedding  $\mathbf{x}_2$  by taking a weighted sum of the value vectors  $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3$ , and  $\mathbf{v}_4$  using the attention weights from the previous step:

$$\mathbf{g}_2 \approx 0.9933 \cdot \mathbf{v}_1 + 0.0067 \cdot \mathbf{v}_2 + 0 \cdot \mathbf{v}_3 + 0 \cdot \mathbf{v}_4$$

As you can see, the decoder’s output for position 2 depends only on (or, we can say “attends only to”) the inputs at positions 1 and 2, with position 2 having a much stronger influence. This effect comes from the **causal mask**, which restricts the model from attending to future positions when generating an output for a given position. This property is essential for maintaining the **autoregressive** nature of language models, ensuring that predictions for each position rely solely on previous and current inputs, not future ones.

While in our example the token primarily attends to itself, this is not a universal pattern. Attention distributions vary based on context, meaning, and token relationships. A token may attend strongly to other tokens that provide relevant semantic or syntactic information, depending on the sentence structure, learned parameters, and specific transformer layer.

The vectors  $\mathbf{q}_t$ ,  $\mathbf{k}_t$ , and  $\mathbf{v}_t$  can be interpreted as follows: each input position (token or embedding) seeks information about other positions. For example, a token like “I” might look for a name in another position, allowing the model to process “I” and the name in a similar way. To enable this, each position  $t$  is assigned a query  $\mathbf{q}_t$ .

The self-attention mechanism calculates a **dot-product** between  $\mathbf{q}_t$  and every key  $\mathbf{k}_p$  across all positions  $p$ . A larger dot-product indicates greater similarity between the vectors. If position  $p$ ’s key  $\mathbf{k}_p$  aligns closely with position  $t$ ’s query  $\mathbf{q}_t$ , then position  $p$ ’s value  $\mathbf{v}_p$  contributes more significantly to the final result.

The concept of attention emerged before the Transformer. In 2014, Dzmitry Bahdanau, while studying under Yoshua Bengio, addressed a fundamental challenge in machine translation: enabling an RNN to focus on the most relevant parts of a sentence. Drawing from his own experience learning English—where he moved his focus between different parts of the text—Bahdanau developed a mechanism<sup>8</sup> for the RNN to “decide” which input words were most important at each translation step. This mechanism, which Bengio termed attention, became a cornerstone of modern neural networks.

---

<sup>8</sup> Bahdanau et al., “Neural Machine Translation by Jointly Learning to Align and Translate,” arXiv preprint, 2014.

The process used to calculate  $\mathbf{g}_2$  is repeated for each position in the input sequence, resulting in a set of output vectors:  $\mathbf{g}_1, \mathbf{g}_2, \mathbf{g}_3$ , and  $\mathbf{g}_4$ . Each position has its own causal mask, so when calculating  $\mathbf{g}_1, \mathbf{g}_3$ , and  $\mathbf{g}_4$ , a different causal mask is applied for each position. The full causal mask for all positions is shown below:

$$\mathbf{M} \stackrel{\text{def}}{=} \begin{bmatrix} 0 & -\infty & -\infty & -\infty \\ 0 & 0 & -\infty & -\infty \\ 0 & 0 & 0 & -\infty \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

As you can see, the first token attends only to itself, the second to itself and the first, the third to itself and the first two, and the last to itself and all preceding tokens.

The general formula for computing attention for all positions is:

$$\mathbf{G} = \text{attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \stackrel{\text{def}}{=} \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} + \mathbf{M}\right)\mathbf{V}$$

Here,  $\mathbf{Q}$  and  $\mathbf{V}$  are  $L \times d_k$  query and value matrices.  $\mathbf{K}^T$  is the  $d_k \times L$  transposed key matrix.  $d_k$  is the dimensionality of the key, query, and value vectors, and  $L$  is the sequence length.

While we computed the attention scores explicitly for  $\mathbf{x}_2$  earlier, the matrix multiplication  $\mathbf{Q}\mathbf{K}^T$  calculates the scores for all positions at once. This method makes the process much faster.

This completes the definition of self-attention.

### 4.3. Position-Wise Multilayer Perceptron

After the masked self-attention layer, each output vector  $\mathbf{g}_t$  is individually processed by a **multilayer perceptron** (MLP). The MLP applies a sequence of additional transformations:

$$\mathbf{z}_t = \mathbf{W}_2(\text{ReLU}(\mathbf{W}_1\mathbf{g}_t + \mathbf{b}_1)) + \mathbf{b}_2$$

Here,  $\mathbf{W}_1$ ,  $\mathbf{W}_2$ ,  $\mathbf{b}_1$ , and  $\mathbf{b}_2$  are learned parameters. The resulting vector  $\mathbf{z}_t$  is then either passed to the next decoder block or, if it's the final decoder block, used to generate the output vector.

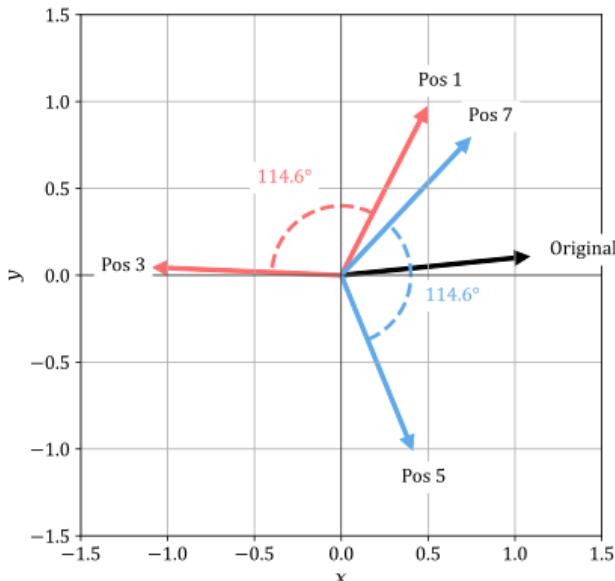
This component is a position-wise multilayer perceptron, which is why I use that term. The literature may refer to it as a feedforward network, dense layer, or fully connected layer, but these names can be misleading. The entire Transformer is a feedforward neural network. Additionally, dense or fully connected layers typically incorporate one weight matrix, one bias vector, and an output non-linearity. The position-wise MLP in a Transformer, however, utilizes two weight matrices, two bias vectors, and omits an output non-linearity.

## 4.4. Rotary Position Embedding

The Transformer architecture, as described so far, does not inherently account for word order. The **causal mask** ensures that each token cannot attend to tokens on its right, but rearranging tokens on the left does not affect the attention weights of a given token. This is unlike RNNs, where hidden states are computed sequentially, each depending on the previous one. Changing word order in RNNs alters the hidden states and, consequently, the output. In contrast, Transformers calculate attention across all tokens at once, without sequential dependency.

To handle word order, Transformers need to incorporate positional information. A widely used method for this is **rotary position embedding (RoPE)**, which applies position-dependent rotations to the query and key vectors in the attention mechanism. One key benefit of RoPE is its ability to generalize effectively to sequences longer than those seen during training. This allows models to be trained on shorter sequences—saving time and computational resources—while still supporting much longer contexts at inference.

RoPE encodes positional information by rotating the query and key vectors. This rotation occurs before the attention computation. Here's a simple illustration of how it works in 2D:



The black arrow labeled “Original” shows a position-less key or query vector in self-attention. RoPE embeds positional information by rotating this vector according to the token’s position.<sup>9</sup> The colored arrows show the resulting rotated vectors for positions 1, 3, 5, and 7.

---

<sup>9</sup> In practice, RoPE operates by rotating pairs of adjacent dimensions within query and key vectors, rather than rotating the entire vectors themselves, as we will explore shortly.

A key property of RoPE is that the angle between any two rotated vectors encodes the distance between their positions in the sequence. For example, the angle between positions 1 and 3 is the same as the angle between positions 5 and 7, since both pairs are two positions apart.

So, how do we rotate vectors? We use matrix multiplication. **Rotation matrices** are widely used in fields like computer graphics to rotate 3D scenes—one of the original purposes of GPUs (the “G” in GPU stands for graphical) before they were applied to neural network training.

In two dimensions, the rotation matrix for an angle  $\theta$  is:

$$\mathbf{R}_\theta = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$$

Let’s rotate the two-dimensional vector  $\mathbf{q} = [2, 1]^\top$ . To do this, we multiply  $\mathbf{q}$  by the rotation matrix  $\mathbf{R}_\theta$ . The result is a new vector, representing  $\mathbf{q}$  rotated counterclockwise by an angle  $\theta$ .

When  $\theta = 45^\circ$  (or  $\pi/4$  radians, since trigonometric functions in PyTorch use radians), we know that  $\cos(\theta) = \sin(\theta) = \frac{\sqrt{2}}{2}$ . Substituting these values, the rotation matrix becomes:

$$\mathbf{R}_{45^\circ} = \begin{bmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix}$$

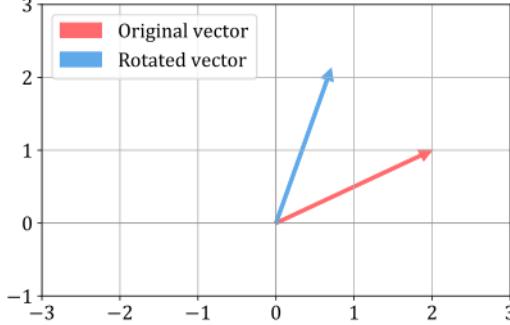
Now, multiplying  $\mathbf{R}_{45^\circ}$  by  $\mathbf{q} = [2, 1]^\top$  gives the rotated vector:

$$\mathbf{R}_{45^\circ} \cdot \mathbf{q} = \begin{bmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

This produces the rotated vector  $\mathbf{q}_{\text{rotated}}$ :

$$\mathbf{q}_{\text{rotated}} = \begin{bmatrix} \frac{\sqrt{2}}{2} \cdot 2 - \frac{\sqrt{2}}{2} \cdot 1 \\ \frac{\sqrt{2}}{2} \cdot 2 + \frac{\sqrt{2}}{2} \cdot 1 \end{bmatrix} = \begin{bmatrix} \frac{\sqrt{2}}{2}(2-1) \\ \frac{\sqrt{2}}{2}(2+1) \end{bmatrix} = \begin{bmatrix} \frac{\sqrt{2}}{2} \cdot 1 \\ \frac{\sqrt{2}}{2} \cdot 3 \end{bmatrix} = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{3\sqrt{2}}{2} \end{bmatrix}$$

The figure below illustrates  $\mathbf{q}$  and its rotated version for  $\theta = 45^\circ$ .



For a position  $t$ , RoPE rotates each pair of dimensions in the query and key vectors defined as:

$$\begin{aligned}\mathbf{q}_t &= [q_t^{(1)}, q_t^{(2)}, \dots, q_t^{(d_q-1)}, q_t^{(d_q)}]^\top \\ \mathbf{k}_t &= [k_t^{(1)}, k_t^{(2)}, \dots, k_t^{(d_k-1)}, k_t^{(d_k)}]^\top\end{aligned}$$

Here,  $d_q$  and  $d_k$  are the (even) dimensionality of the query and key vectors. RoPE rotates pairs of dimensions indexed as  $(2p-1, 2p)$ , where  $p$  ranges from 1 to  $d_q/2$  and represents each pair's index.

To split the dimensions of  $\mathbf{q}_t$  into  $d_q/2$  pairs, we group them as:

$$[q_t^{(1)}, q_t^{(2)}]^\top, [q_t^{(3)}, q_t^{(4)}]^\top, \dots, [q_t^{(d_q-1)}, q_t^{(d_q)}]^\top$$

When we write  $\mathbf{q}_t(p)$ , it represents the pair  $[q_t^{(2p-1)}, q_t^{(2p)}]$ . For example,  $\mathbf{q}_t(3)$  corresponds to:

$$[q_t^{(2\cdot 3-1)}, q_t^{(2\cdot 3)}] = [q_t^{(5)}, q_t^{(6)}]$$

Each pair  $p$  undergoes a rotation based on the token position  $t$  and a **rotation frequency**  $\theta_p$ :

$$\text{RoPE}(\mathbf{q}_t(p)) \stackrel{\text{def}}{=} \begin{bmatrix} \cos(\theta_p t) & -\sin(\theta_p t) \\ \sin(\theta_p t) & \cos(\theta_p t) \end{bmatrix} \begin{bmatrix} q_t^{(2p-1)} \\ q_t^{(2p)} \end{bmatrix}$$

Applying the **matrix-vector multiplication** rule, the rotation results in the following 2D vector:

$$\text{RoPE}(\mathbf{q}_t(p)) = [q_t^{(2p-1)} \cos(\theta_p t) - q_t^{(2p)} \sin(\theta_p t), q_t^{(2p-1)} \sin(\theta_p t) + q_t^{(2p)} \cos(\theta_p t)]^\top,$$

where  $\theta_p$  is the rotation frequency for the  $p^{\text{th}}$  pair. It is defined as:

$$\theta_p = \frac{1}{\Theta^{2(p-1)/d_q}}$$

Here,  $\Theta$  is a constant. Initially set to 10,000, later experiments demonstrated that higher values of  $\Theta$ —such as 500,000 (used in Llama 2 and 3 series of models) or 1,000,000 (in Qwen 2 and 2.5 series)—enable support for larger context sizes (hundreds of thousands of tokens).

The full rotated embedding  $\text{RoPE}(\mathbf{q}_t)$  is constructed by concatenating all the rotated pairs:

$$\text{RoPE}(\mathbf{q}_t) \stackrel{\text{def}}{=} \text{concat}\left[\text{RoPE}(\mathbf{q}_t(1)), \text{RoPE}(\mathbf{q}_t(2)), \dots, \text{RoPE}(\mathbf{q}_t(d_q/2))\right]$$

Note how the rotation frequency  $\theta_p$  decreases quickly for each subsequent pair because of the exponential term in the denominator. This enables RoPE to capture fine-grained local position information in the early dimensions, where rotations are more frequent, and coarse-grained global position information in the later dimensions, where rotations slow down. This combination creates richer positional encoding, allowing the model to differentiate token positions in a sequence more effectively than using a single rotation frequency across all dimensions.

To illustrate the process, consider a 6-dimensional query vector at position  $t$  and  $\theta = 10,000$ :

$$\mathbf{q}_t = [q_t^{(1)}, q_t^{(2)}, q_t^{(3)}, q_t^{(4)}, q_t^{(5)}, q_t^{(6)}]^\top \stackrel{\text{def}}{=} [0.8, 0.6, 0.7, 0.3, 0.5, 0.4]^\top$$

First, we split it into three pairs ( $d_q/2 = 3$ ):

$$\begin{aligned}\mathbf{q}_t(1) &= [q_t^{(1)}, q_t^{(2)}] = [0.8, 0.6]^\top \\ \mathbf{q}_t(2) &= [q_t^{(3)}, q_t^{(4)}] = [0.7, 0.3]^\top \\ \mathbf{q}_t(3) &= [q_t^{(5)}, q_t^{(6)}] = [0.5, 0.4]^\top\end{aligned}$$

Each pair  $p$  undergoes a rotation by angle  $\theta_p t$ , where:

$$\theta_p = \frac{1}{10000^{2(p-1)/d_q}}$$

Let the position  $t$  be 100. First, we calculate the rotation angles for each pair (in radians):

$$\begin{aligned}\theta_1 &= \frac{1}{10000^{2(1-1)/6}} = \frac{1}{10000^{0/6}} = 1.0000, \quad \text{therefore: } \theta_1 t = 100.00 \\ \theta_2 &= \frac{1}{10000^{2(2-1)/6}} = \frac{1}{10000^{2/6}} \approx 0.0464, \quad \text{therefore: } \theta_2 t = 4.64 \\ \theta_3 &= \frac{1}{10000^{2(3-1)/6}} = \frac{1}{10000^{4/6}} \approx 0.0022, \quad \text{therefore: } \theta_3 t = 0.22\end{aligned}$$

The rotated pair 1 is:

$$\text{RoPE}(\mathbf{q}_{100}(1)) = \begin{bmatrix} \cos(100) & -\sin(100) \\ \sin(100) & \cos(100) \end{bmatrix} \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix} \approx \begin{bmatrix} 0.86 & 0.51 \\ -0.51 & 0.86 \end{bmatrix} \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix} = [0.99, 0.11]^\top$$

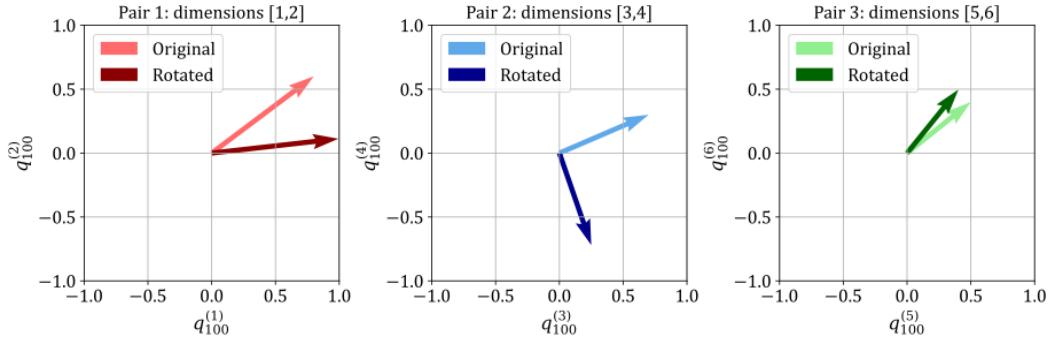
The rotated pair 2 is:

$$\text{RoPE}(\mathbf{q}_{100}(2)) = \begin{bmatrix} \cos(4.64) & -\sin(4.64) \\ \sin(4.64) & \cos(4.64) \end{bmatrix} \begin{bmatrix} 0.7 \\ 0.3 \end{bmatrix} \approx \begin{bmatrix} -0.07 & 1.00 \\ -1.00 & -0.07 \end{bmatrix} \begin{bmatrix} 0.7 \\ 0.3 \end{bmatrix} = [0.25, -0.72]^\top$$

The rotated pair 3 is:

$$\text{RoPE}(\mathbf{q}_{100}(3)) = \begin{bmatrix} \cos(0.22) & -\sin(0.22) \\ \sin(0.22) & \cos(0.22) \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.4 \end{bmatrix} \approx \begin{bmatrix} 0.98 & -0.21 \\ 0.21 & 0.98 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.4 \end{bmatrix} = [0.40, 0.50]^\top$$

These is what the original and rotated pairs look like when plotted:



The final RoPE-encoded vector is the concatenation of these pairs:

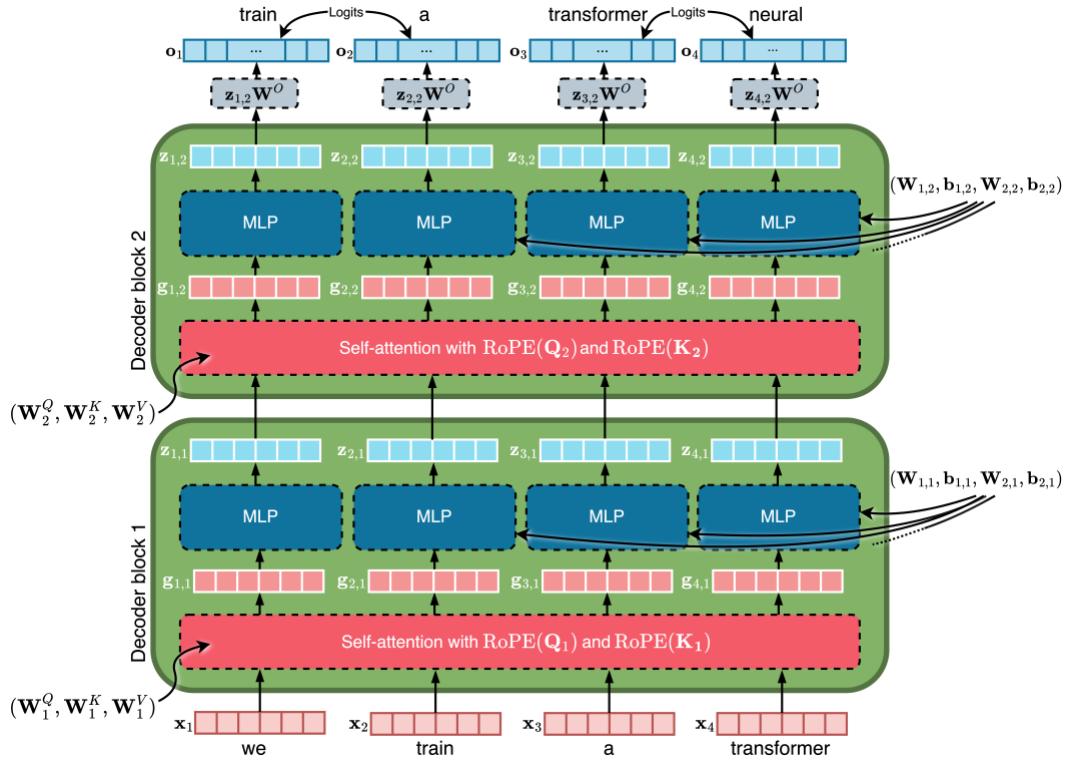
$$\text{RoPE}(\mathbf{q}_{100}) \approx [0.99, 0.11, 0.25, -0.72, 0.40, 0.50]^\top$$

The math for  $\text{RoPE}(\mathbf{k}_t)$  is the same as for  $\text{RoPE}(\mathbf{q}_t)$ . In each decoder block, RoPE is applied to each row of the query ( $\mathbf{Q}$ ) and key ( $\mathbf{K}$ ) matrices within the self-attention mechanism.

Value vectors only provide the information that is selected and combined after the attention weights are determined. Since the positional relationships are already captured in the query-key alignment, value vectors don't need their own rotary embeddings. In other words, the value vectors simply "deliver" the content once the positional-aware attention has identified where to look.

Recall that  $\mathbf{Q}$  and  $\mathbf{K}$  are generated by multiplying the decoder block inputs by weight matrices  $\mathbf{W}^Q$  and  $\mathbf{W}^K$ , as illustrated in Figure 4.1. RoPE is applied immediately after obtaining  $\mathbf{Q}$  and  $\mathbf{K}$ , and before the attention scores are calculated.

Applying RoPE in all decoder blocks helps retain positional information throughout the network's depth. The illustration below shows two decoder blocks and where RoPE is applied:



In this illustration, the outputs of the second decoder block are used to compute logits for each position. This is achieved by multiplying the outputs of the final decoder block by a matrix of shape (embedding dimensionality, vocabulary size), which is shared across all positions. We will implement the decoder model in Python soon, where this detail will become clearer.

The self-attention mechanism we've described would work as is. However, transformers typically employ an enhanced version called **multi-head attention**. This allows the model to focus on multiple aspects of information simultaneously. For example, one attention head might capture syntactic relationships, another might emphasize semantic similarities, and a third could detect long-range dependencies between tokens.

## 4.5. Multi-Head Attention

Once you understand self-attention, understanding multi-head attention is relatively straightforward. For each head, from 1 to  $H$ , there is a separate triplet of attention matrices:

$$\{(\mathbf{W}_h^Q, \mathbf{W}_h^K, \mathbf{W}_h^V)\}_{h \in 1, \dots, H}$$

Each triplet is applied to the input vectors  $\mathbf{x}_1, \dots, \mathbf{x}_4$ , producing  $H$  matrices  $\mathbf{G}_h$ . For each head  $h$ , this gives four vectors  $\mathbf{g}_{h,1}, \dots, \mathbf{g}_{h,4}$ , as shown in Figure 4.2 for three heads ( $H = 3$ ):

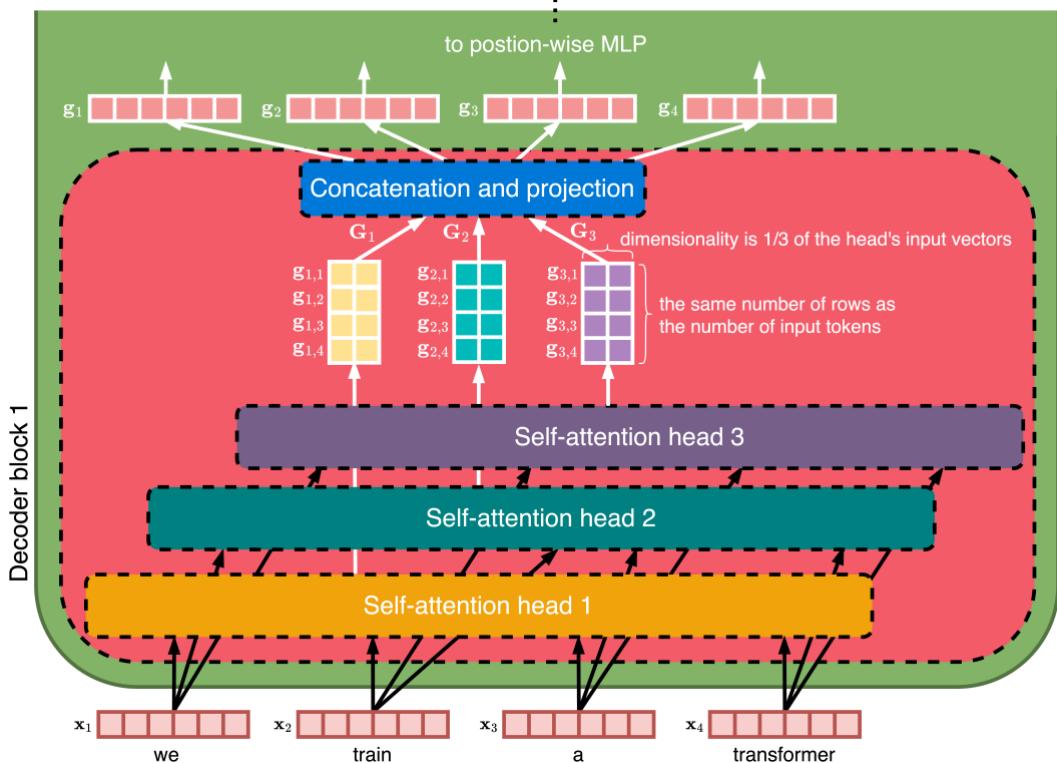
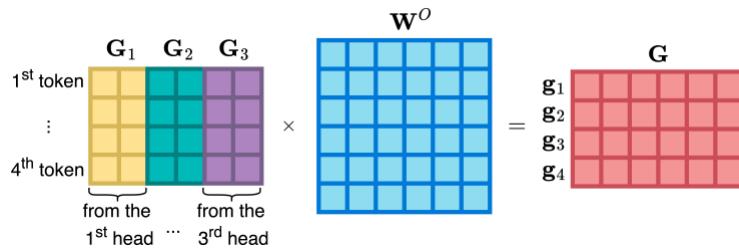


Figure 4.2: 3-head self-attention.

As you can see, the multi-head self-attention mechanism processes an input sequence through multiple self-attention “heads.” For instance, with 3 heads, each head calculates self-attention scores for the input tokens independently. RoPE is applied separately in each head.

All input tokens  $x_1, \dots, x_4$  are processed by all three heads, producing output matrices  $G_1, G_2$ , and  $G_3$ . Each matrix  $G_h$  has as many rows as there are input tokens, meaning each head generates an embedding for every token. The embedding dimensionality of each  $G_h$  is reduced to one-third of the total embedding dimensionality. As a result, each head outputs lower-dimensional embeddings compared to the original embedding size.

The outputs from the three heads are concatenated along the embedding dimension in the **concatenation and projection layer**, creating a single matrix that integrates information from all heads. This matrix is then transformed by the **projection matrix  $W^O$** , resulting in the final output matrix  $G$ . This output is passed to the position-wise MLP:



Concatenating the matrices  $G_1$ ,  $G_2$ , and  $G_3$  restores the original embedding dimensionality (e.g., 6 in this case). However, applying the trainable parameter matrix  $W^O$  enables the model to combine the heads' information more effectively than mere concatenation.

The original Transformer paper tested various numbers of attention heads and found the optimal range to be around 8–16 heads. In contrast, modern large language models often use 16 to 64 heads.

At this stage, the reader understands the Transformer model architecture at a high level. Two key technical details remain to explore: layer normalization and residual connections, both essential components that enable the Transformer's effectiveness. Let's begin with residual connections.

## 4.6. Residual Connection

**Residual connections**, or **skip connections**, form a key component of the Transformer architecture. They address vanishing gradients in deep neural networks, making it possible to train substantially deeper models.

A network containing more than two layers is called a **deep neural network**. The process of training these models is known as **deep learning**. Prior to the development of **ReLU** and residual connections, training deep networks posed significant challenges due to the **vanishing gradient problem**. Let's examine this phenomenon.

Remember that in the gradient descent algorithm, we calculate partial derivatives for all parameters and update them by taking a small step in the direction opposite to the gradient. As networks grow deeper, this step becomes progressively smaller in the earlier layers (those closer to the input), resulting in minimal (vanishing) parameter updates in these layers. Residual connections strengthen these updates by creating pathways for the gradient to "bypass" certain layers, hence the term skip connections.

To understand the vanishing gradient problem more thoroughly, let's analyze a 3-layer neural network expressed as a **composite function**:

$$f(x) = f_3(f_2(f_1(x))),$$

where  $f_1$  represents the first layer,  $f_2$  represents the second layer, and  $f_3$  represents the third (output) layer. Let these functions be defined as follows:

$$\begin{aligned} z &\leftarrow f_1(x) \stackrel{\text{def}}{=} w_1 x + b_1 \\ r &\leftarrow f_2(z) \stackrel{\text{def}}{=} w_2 z + b_2 \\ y &\leftarrow f_3(r) \stackrel{\text{def}}{=} w_3 r + b_3 \end{aligned}$$

Here,  $w_l$  and  $b_l$  are scalar weights and biases for each layer  $l \in \{1, 2, 3\}$  and the notation  $z \leftarrow f_1(x)$  means  $f_1(x)$  takes  $x$  as input and returns  $z$ .

Let's define the loss function  $L$  in terms of the network output  $f(x)$  and the true label  $y$  as  $L(f(x), y)$ . The gradient of the loss  $L$  with respect to  $w_1$ , denoted as  $\frac{\partial L}{\partial w_1}$ , is given by:

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial w_1} = \frac{\partial L}{\partial f_3} \cdot \frac{\partial f_3}{\partial f_2} \cdot \frac{\partial f_2}{\partial f_1} \cdot \frac{\partial f_1}{\partial w_1},$$

where:

$$\frac{\partial f_3}{\partial f_2} = w_3, \quad \frac{\partial f_2}{\partial f_1} = w_2, \quad \frac{\partial f_1}{\partial w_1} = x$$

So, we can write:

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f_3} \cdot w_3 \cdot w_2 \cdot x$$

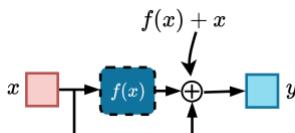
The vanishing gradient problem occurs when weights like  $w_2$  and  $w_3$  are small (less than 1). When multiplied together, they produce even smaller values, causing the gradient for earlier weights such as  $w_1$  to approach zero. This issue becomes particularly severe in deep networks with many layers.

Take large language models as an example. These networks often include 32 or more decoder blocks. To simplify, assume all blocks are fully connected layers. If the average weight value is around 0.5, the gradient for the input layer parameters becomes  $0.5^{32} \approx 0.0000000002$ . This is extremely small. After multiplying by the learning rate, updates to the early layers are negligible. As a result, the network stops learning effectively.

Residual connections offer a solution to the vanishing gradient problem by creating shortcuts in the gradient computation path. The basic idea is simple: instead of passing only the output of a layer to the next one, the layer's input is added to its output. Mathematically, this is written as:

$$y = f(x) + x,$$

where  $x$  is the input,  $f(x)$  is the layer's computed function, and  $y$  is the output. This addition forms the residual connection. Graphically, it looks like this:



In this illustration, the input  $x$  is processed both through the layer (represented as  $f(x)$ ) and added directly to the layer's output.

Now let's introduce residual connections into our 3-layer network, where inputs and outputs are scalars. We'll see how this changes gradient computation and mitigates the vanishing gradient issue. Starting with the original network  $f(x) = f_3(f_2(f_1(x)))$ , let's add residual connections to layers 2 and 3:

$$\begin{aligned} z &\leftarrow f_1(x) \stackrel{\text{def}}{=} w_1x + b_1 \\ r &\leftarrow f_2(z) \stackrel{\text{def}}{=} w_2z + b_2 + z \\ y &\leftarrow f_3(r) \stackrel{\text{def}}{=} w_3r + b_3 + r \end{aligned}$$

Our composite function becomes:

$$f(x) = w_3[w_2(w_1x + b_1) + b_2 + w_1x + b_1] + b_3 + w_2(w_1x + b_1) + b_2 + w_1x + b_1$$

Now, let's calculate the gradient of the loss  $L$  with respect to  $w_1$ :

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial w_1}$$

Expanding  $\frac{\partial f}{\partial w_1}$ :

$$\begin{aligned} \frac{\partial f}{\partial w_1} &= \frac{\partial}{\partial w_1} [(w_3(w_2(w_1x + b_1) + b_2 + (w_1x + b_1)) + b_3) + (w_2(w_1x + b_1) + b_2 + (w_1x + b_1))] \\ &= (w_3w_2 + w_3 + w_2 + 1) \cdot x \end{aligned}$$

Therefore, the full gradient is:

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f} \cdot (w_3w_2 + w_3 + w_2 + 1) \cdot x$$

Comparing this to our original gradient without residual connections:

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f} \cdot w_3 \cdot w_2 \cdot x$$

We observe that residual connections introduce three additional terms:  $w_3$ ,  $w_2$ , and 1. This guarantees that the gradient will not vanish completely, even when  $w_2$  and  $w_3$  are small, due to the added constant term 1.

For example, if  $w_2 = w_3 = 0.5$  as in the previous case:

- **Without residual connections:**  $0.5 \cdot 0.5 = 0.25$
- **With residual connections:**  $0.5 \cdot 0.5 + 0.5 + 0.5 + 1 = 2.25$

The illustration below depicts an encoding block with residual connections:



As shown, each decoder block includes two residual connections. The layers are now named like Python objects, which we will implement shortly. Additionally, two RMSNorm layers have been added. Let's discuss their purpose.

## 4.7. Root Mean Square Normalization

The RMSNorm layer applies **root mean square normalization** to the input vector. This operation takes place just before the vector enters the self-attention layer and the position-wise MLP. Let's illustrate this with a three-dimensional vector.

Suppose we have a vector  $\mathbf{x} = [x^{(1)}, x^{(2)}, x^{(3)}]^T$ . To apply RMS normalization, we first calculate the **root mean square (RMS)** of the vector:

$$\text{RMS}(\mathbf{x}) = \sqrt{\frac{1}{3} \sum_{i=1}^3 (x^{(i)})^2} = \sqrt{\frac{1}{3} [(x^{(1)})^2 + (x^{(2)})^2 + (x^{(3)})^2]}$$

Then, we normalize the vector by dividing each component by the RMS value to obtain  $\tilde{\mathbf{x}}$ :

$$\tilde{\mathbf{x}} = \frac{\mathbf{x}}{\text{RMS}(\mathbf{x})} = \left[ \frac{x^{(1)}}{\text{RMS}(\mathbf{x})}, \frac{x^{(2)}}{\text{RMS}(\mathbf{x})}, \frac{x^{(3)}}{\text{RMS}(\mathbf{x})} \right]^\top$$

Finally, we apply the scale factor  $\gamma$  to each dimension of  $\tilde{\mathbf{x}}$ :

$$\bar{\mathbf{x}} = \text{RMSNorm}(\mathbf{x}) \stackrel{\text{def}}{=} \boldsymbol{\gamma} \odot \tilde{\mathbf{x}} = [\gamma^{(1)} \tilde{x}^{(1)}, \gamma^{(2)} \tilde{x}^{(2)}, \gamma^{(3)} \tilde{x}^{(3)}]^\top,$$

where  $\odot$  denotes the **element-wise product**. The vector  $\boldsymbol{\gamma}$  is a trainable parameter, and each RMSNorm layer has its own independent  $\boldsymbol{\gamma}$ .

The primary purpose of RMSNorm is to stabilize training by keeping the scale of the input to each layer consistent. This improves numerical stability, helping to prevent gradient updates that are excessively large or small.

Now that we've covered the key components of the Transformer architecture, let's summarize how a decoder block processes its input:

1. The input embeddings  $\mathbf{x}_t$  first go through RMS normalization.
2. The normalized embeddings  $\tilde{\mathbf{x}}_t$  are processed by the multi-head self-attention mechanism, with RoPE applied to key and query vectors.
3. The self-attention output  $\mathbf{g}_t$  is added to the original input  $\mathbf{x}_t$  (residual connection).
4. This sum,  $\hat{\mathbf{g}}_t$ , undergoes RMS normalization again.
5. The normalized sum  $\bar{\mathbf{g}}_t$  is passed through the multilayer perceptron.
6. The perceptron output  $\mathbf{z}_t$  is added to the pre-RMS-normalization vector  $\hat{\mathbf{g}}_t$  (another residual connection).
7. The result,  $\hat{\mathbf{z}}_t$ , is the output of the decoder block, serving as input for the next block (or the final output layer if it's the last block).

This sequence is repeated for each decoder block in the Transformer.

## 4.8. Key-Value Caching

During training, the decoder can process all positions in parallel because at each layer it computes the query, key, and value matrices,  $\mathbf{Q} = \mathbf{X}\mathbf{W}^Q$ ,  $\mathbf{K} = \mathbf{X}\mathbf{W}^K$ , and  $\mathbf{V} = \mathbf{X}\mathbf{W}^V$ , for the entire sequence  $\mathbf{X}$ . However, in autoregressive (left-to-right) generation, tokens must be generated one at a time. Normally, each time we generate a new token, we would have to:

1. Recalculate the key and value matrices for all previous tokens.
2. Merge these with the new token's key and value vectors to compute self-attention.

**Key-value caching** prevents this redundant recomputation. Since  $\mathbf{W}^K$  and  $\mathbf{W}^V$  stay fixed after training, the key and value vectors for earlier tokens remain unchanged during decoding. This allows us to store (“cache”) these vectors after computing them once. For each new token:

- We only compute its key and value vectors using  $\mathbf{W}^K$  and  $\mathbf{W}^V$ .
- We then append these vectors to the cached key-value pairs for self-attention.

This approach ensures that the rest of the sequence does not need to be reprocessed, which reduces computation significantly—especially for long sequences. For each decoder block, the cached keys and values are stored per attention head, with shapes  $(L \times d_h)$  for both matrices (where  $L$  increases by one with each new token and  $d_h$  is the dimensionality of the query, key, and value vectors for this attention head). Consequently, for a model with  $H$  attention heads, the combined key and value caches in each decoder block have shapes  $(H \times L \times d_h)$ .

Now that we understand how the Transformer operates, we’re ready to start coding.

## 4.9. Transformer in Python

Let’s begin implementing the decoder in Python by defining the `AttentionHead` class:

```
class AttentionHead(nn.Module):
    def __init__(self, emb_dim, d_h):
        super().__init__()
        self.W_Q = nn.Parameter(torch.empty(emb_dim, d_h))
        self.W_K = nn.Parameter(torch.empty(emb_dim, d_h))
        self.W_V = nn.Parameter(torch.empty(emb_dim, d_h))
        self.d_h = d_h

    def forward(self, x, mask):
        Q = x @ self.W_Q ❶
        K = x @ self.W_K
        V = x @ self.W_V ❷

        Q, K = rope(Q), rope(K) ❸

        scores = Q @ K.transpose(-2, -1) / math.sqrt(self.d_h) ❹
        masked_scores = scores.masked_fill(mask == 0, float("-inf")) ❺
        attention_weights = torch.softmax(masked_scores, dim=-1) ❻
        return attention_weights @ V ❼
```

This class implements a single attention head in the multi-head attention mechanism. In the constructor, we initialize three trainable weight matrices: the query matrix  $\mathbf{W}_Q$ , the key matrix  $\mathbf{W}_K$ , and the value matrix  $\mathbf{W}_V$ . Each of these is a `Parameter` tensor of shape  $(\text{emb\_dim}, d_h)$ , where `emb_dim` is the input embedding dimension and `d_h` is the dimensionality of the query, key, and value vectors for this attention head.

In the `forward` method:

- Lines ❶ and ❷ compute the query, key, and value matrices by multiplying the input vector  $x$  with the respective weight matrices. Given that  $x$  has shape  $(\text{batch\_size}, \text{seq\_len}, \text{emb\_dim})$ ,  $Q$ ,  $K$ , and  $V$  each have shape  $(\text{batch\_size}, \text{seq\_len}, d_h)$ .
- Line ❸ applies the rotary positional encoding to  $Q$  and  $K$ . After the query and key vectors are rotated, line ❹ computes the attention scores. Here's a breakdown:
  - $K.\text{transpose}(-2, -1)$  swaps the last two dimensions of  $K$ . If  $K$  has shape  $(\text{batch\_size}, \text{seq\_len}, d_h)$ , transposing it results in  $(\text{batch\_size}, d_h, \text{seq\_len})$ . This prepares  $K$  for matrix multiplication with  $Q$ .
  - $Q @ K.\text{transpose}(-2, -1)$  performs batch matrix multiplication, resulting in a tensor of attention scores of shape  $(\text{batch\_size}, \text{seq\_len}, \text{seq\_len})$ .
  - As mentioned in Section 4.2, we divide by  $\sqrt{d_h}$  for numerical stability.

When the matrix multiplication operator `@` is applied to tensors with more than two dimensions, PyTorch uses **broadcasting**. This technique handles dimensions that aren't directly compatible with the `@` operator, which is normally defined only for two-dimensional tensors (matrices). In this case, PyTorch treats the first dimension as the batch dimension, performing the matrix multiplication separately for each example in the batch. This process is known as **batch matrix multiplication**.

- Line ❺ applies the causal mask. The `mask` tensor has the shape  $(\text{seq\_len}, \text{seq\_len})$  and contains 0s and 1s. The `masked_fill` function replaces all cells in the input matrix where `mask == 0` with negative infinity. This prevents attention to future tokens. Since the `mask` lacks the batch dimension while `scores` includes it, PyTorch uses broadcasting to apply the `mask` to the `scores` of each sequence in the batch.
- Line ❻ applies softmax to the `scores` along the last dimension, turning them into attention weights. Then, line ❼ computes the output by multiplying these attention weights with  $V$ . The resulting output has the shape  $(\text{batch\_size}, \text{seq\_len}, d_h)$ .

Given the attention head class, we can now define the `MultiHeadAttention` class:

```
class MultiHeadAttention(nn.Module):
    def __init__(self, emb_dim, num_heads):
        super().__init__()
        d_h = emb_dim // num_heads ❶
        self.heads = nn.ModuleList([
            AttentionHead(emb_dim, d_h)
            for _ in range(num_heads)
        ]) ❷
        self.W_O = nn.Parameter(torch.empty(emb_dim, emb_dim)) ❸

    def forward(self, x, mask):
        head_outputs = [head(x, mask) for head in self.heads] ❹
```

```

x = torch.cat(head_outputs, dim=-1) ❸
return x @ self.W_O ❹

```

In the constructor:

- Line ❶ calculates  $d_h$ , the dimensionality of each attention head, by dividing the model's embedding dimensionality  $\text{emb\_dim}$  by the number of heads.
- Line ❷ creates a `ModuleList` containing  $\text{num\_heads}$  instances of `AttentionHead`. Each head takes the input dimensionality  $\text{emb\_dim}$  and outputs a vector of size  $d_h$ .
- Line ❸ initializes  $W_O$ , a learnable **projection matrix** with shape  $(\text{emb\_dim}, \text{emb\_dim})$  to combine the outputs from all attention heads.

In the `forward` method:

- Line ❹ applies each attention head to the input  $x$  of shape  $(\text{batch\_size}, \text{seq\_len}, \text{emb\_dim})$ . Each head's output has shape  $(\text{batch\_size}, \text{seq\_len}, d_h)$ .
- Line ❺ concatenates all heads' outputs along the last dimension. The resulting  $x$  has shape  $(\text{batch\_size}, \text{seq\_len}, \text{emb\_dim})$  since  $\text{num\_heads} * d_h = \text{emb\_dim}$ .
- Line ❻ multiplies the concatenated output by the projection matrix  $W_O$ . The output has the same shape as input.

Now that we have multi-head attention, the last piece needed for the decoder block is the position-wise multilayer perceptron. Let's define it:

```

class MLP(nn.Module):
    def __init__(self, emb_dim):
        super().__init__()
        self.W_1 = nn.Parameter(torch.empty(emb_dim, emb_dim * 4))
        self.B_1 = nn.Parameter(torch.empty(emb_dim * 4))
        self.W_2 = nn.Parameter(torch.empty(emb_dim * 4, emb_dim))
        self.B_2 = nn.Parameter(torch.empty(emb_dim))

    def forward(self, x):
        x = x @ self.W_1 + self.B_1 ❶
        x = torch.relu(x) ❷
        x = x @ self.W_2 + self.B_2 ❸
        return x

```

In the constructor, we initialize learnable weights and biases.

In the `forward` method:

- Line ❶ multiplies the input  $x$  by the weight matrix  $W_1$  and adds the bias vector  $B_1$ . The input has shape  $(\text{batch\_size}, \text{seq\_len}, \text{emb\_dim})$ , so the result has shape  $(\text{batch\_size}, \text{seq\_len}, \text{emb\_dim} * 4)$ .
- Line ❷ applies the **ReLU** activation function element-wise, adding non-linearity.

- Line ❸ multiplies the result by the second weight matrix  $W_2$  and adds the bias vector  $B_2$ , reducing the dimensionality back to  $(batch\_size, seq\_len, emb\_dim)$ .

The first linear transformation expands to 4 times the embedding dimensionality ( $emb\_dim * 4$ ) to provide the network with greater capacity for learning complex patterns and relationships between variables. The 4x factor balances expressiveness and efficiency.

After expanding the dimensionality, it's compressed back to the original embedding dimensionality ( $emb\_dim$ ). This ensures compatibility with residual connections, which require matching dimensionalities. Empirical results support this expand-and-compress approach as an effective trade-off between computational cost and performance.

With all components defined, we're ready to set up the complete decoder block:

```
class DecoderBlock(nn.Module):
    def __init__(self, emb_dim, num_heads):
        super().__init__()
        self.norm1 = RMSNorm(emb_dim)
        self.attn = MultiHeadAttention(emb_dim, num_heads)
        self.norm2 = RMSNorm(emb_dim)
        self.mlp = MLP(emb_dim)

    def forward(self, x, mask):
        attn_out = self.attn(self.norm1(x), mask) ❶
        x = x + attn_out ❷
        mlp_out = self.mlp(self.norm2(x)) ❸
        x = x + mlp_out ❹
        return x
```

The `DecoderBlock` class represents a single decoder block in a Transformer model. In the constructor, we set up the necessary layers: two `RMSNorm` layers, a `MultiHeadAttention` instance (configured with the embedding dimensionality and number of heads), and an `MLP` layer.

In the `forward` method:

- Line ❶ applies `RMSNorm` to the input  $x$ , which has shape  $(batch\_size, seq\_len, emb\_dim)$ . The output of `RMSNorm` keeps this shape. This normalized tensor is then passed to the multi-head attention layer, which outputs a tensor of the same shape.
- Line ❷ adds a **residual connection** by combining the attention output `attn_out` with the original input  $x$ . The shape doesn't change.
- Line ❸ applies the second `RMSNorm` to the result from the residual connection, retaining the same shape. This normalized tensor is then passed through the `MLP`, which outputs another tensor with shape  $(batch\_size, seq\_len, emb\_dim)$ .
- Line ❹ adds a second residual connection, combining `mlp_out` with its unnormalized input. The decoder block's final output shape is  $(batch\_size, seq\_len, emb\_dim)$ , ready for the next decoder block or the final output layer.

With the decoder block defined, we can now build the decoder transformer language model by stacking multiple decoder blocks sequentially:

```
class DecoderLanguageModel(nn.Module):
    def __init__(self, vocab_size, emb_dim, num_heads, num_blocks, pad_idx):
        super().__init__()
        self.embedding = nn.Embedding(vocab_size, emb_dim, padding_idx=pad_idx)
        self.layers = nn.ModuleList([
            DecoderBlock(emb_dim, num_heads) for _ in range(num_blocks)
        ])
        self.output = nn.Parameter(torch.rand(emb_dim, vocab_size))

    def forward(self, x):
        x = self.embedding(x)
        _, seq_len, _ = x.shape
        mask = torch.tril(torch.ones(seq_len, seq_len, device=x.device))
        for layer in self.layers:
            x = layer(x, mask)
        return x @ self.output
```

In the constructor of the `DecoderLanguageModel` class:

- Line ❶ creates an embedding layer that converts input token indices to dense vectors. The `padding_idx` specifies the ID of the padding token, ensuring that padding tokens are mapped to zero vectors.
- Line ❷ creates a `ModuleList` with `num_blocks` `DecoderBlock` instances, forming the stack of decoder layers.
- Line ❸ defines a matrix to project the last decoder block's output to logits over the vocabulary, enabling next token prediction.

In the `forward` method:

- Line ❹ converts the input token indices to embeddings. The input tensor `x` has shape `(batch_size, seq_len)`; the output has shape `(batch_size, seq_len, emb_dim)`.
- Line ❺ creates the **causal mask**.
- Line ❻ applies each decoder block to the input tensor `x` with shape `(batch_size, seq_len, emb_dim)`, producing an output tensor of the same shape. Each block refines the sequence and passes it to the next until the final block.
- Line ❼ projects the output of the final decoder block to vocabulary-sized logits by multiplying it with the `self.output` matrix, which has shape `(emb_dim, vocab_size)`.

After this batched matrix multiplication, the final output has shape `(batch_size, seq_len, vocab_size)`, providing scores for each token in the vocabulary at each position in the input sequence. This output can then be used to generate the model's predictions as we will discuss in the next chapter.

The training loop for `DecoderLanguageModel` is the same as for the RNN (Section 3.6), so it is not repeated here for brevity. Implementations of RMSNorm and RoPE are also skipped. Training data is prepared just like for the RNN: the target sequence is offset by one position relative to the input sequence, as described in Section 3.7. The complete code for training the decoder language model is available in the [thelmbbook.com/nb/4.1](http://thelmbbook.com/nb/4.1) notebook.

In the notebook, I used these hyperparameter values: `emb_dim = 128`, `num_heads = 8`, `num_blocks = 2`, `batch_size = 128`, `learning_rate = 0.001`, `num_epochs = 1`, and `context_size = 30`. With these settings, the model achieved a perplexity of 55.19, improving on the RNN's 72.23. This is a good result given the comparable number of trainable parameters (8,621,963 for the Transformer vs. 8,292,619 for the RNN). The real strengths of transformers, however, become apparent at larger scales of model size, context length, and training data. Reproducing experiments at such scales in this book is impractical.

Let's look at some continuations of the prompt "The President" generated by the decoder model at later training steps:

```
The President has been in the process of a new deal to make a decision on the issue .
```

```
The President 's office said the government had `` no intention of making any mistakes '' .
```

```
The President of the United States has been a key figure for the first time in the past ## years .
```

The "#" characters in the training data represent individual digits. For example, "##" likely represents the number of years.

---

If you've made it this far, well done! You now understand the mechanics of language models. But understanding the mechanics alone won't help you fully appreciate what modern language models are capable of. To truly understand, you need to work with one.

In the next chapter, we'll explore large language models (LLMs). We'll discuss why they're called *large* and what's so special about the size. Then, we'll cover how to finetune an existing LLM for practical tasks like question answering and document classification, as well as how to use LLMs to address a variety of real-world problems.

# Chapter 5. Large Language Model

Large language models have transformed NLP through their remarkable capabilities in text generation, translation, and question-answering. But how can a model trained solely to predict the next word achieve these results? The answer lies in two factors: scale and supervised finetuning.

## 5.1. Why Larger Is Better

LLMs are built with a large number of parameters, large context windows, and trained on large corpora backed by substantial computational resources. This scale enables them to learn complex language patterns and even memorize information.

Creating a **chat LM**, capable of handling dialogue and following complex instructions, involves two stages. The first stage is **pretraining** on a massive dataset, often containing trillions of tokens. In this phase, the model learns to predict the next token based on context—similar to what we did with the RNN and decoder models, but at a vastly larger scale.

With more parameters and extended context windows, the model aims to “understand” the context as deeply as possible to improve the next token prediction and minimize the **cross-entropy** loss. For example, consider this context:

The CRISPR-Cas9 technique has revolutionized genetic engineering by enabling precise modifications to DNA sequences. The process uses a guide RNA to direct the Cas9 enzyme to a specific location in the genome. Once positioned, Cas9 acts like molecular scissors, cutting the DNA strand. This cut activates the cell's natural repair mechanisms, which scientists can exploit to

To accurately predict the next token, the model must know:

- 1) about CRISPR-Cas9 and its components, such as guide RNA and Cas9 enzyme,
- 2) how CRISPR-Cas9 works—locating specific DNA sequences and cutting DNA,
- 3) about cellular repair mechanisms, and
- 4) how these mechanisms enable gene editing.

A well-trained LLM might suggest continuations like “insert new genetic material” or “delete unwanted genes.” Choosing “insert” or “delete” over vague terms like “change” or “fix” requires encoding the context into embedding vectors that reflect a deeper understanding of the gene-editing process, rather than relying on surface-level patterns as count-based models do.

It’s intuitive to think that if words and paragraphs can be represented by dense embedding vectors, then entire documents or complex explanations could theoretically be represented this way too. However, before LLMs were discovered, NLP researchers believed embeddings could only represent basic concepts like “animal,” “building,” “economy,” “technology,” “verb,” or “noun.” This belief is evident in the conclusion of one of the most influential papers of the 2010s, which detailed the training of a state-of-the-art language model at that time:

*"As with all text generated by language models, the sample does not make sense beyond the level of short phrases. The realism could perhaps be improved with a larger network and/or more data. However, it seems futile to expect meaningful language from a machine that has never been exposed to the sensory world to which language refers."* (Alex Graves, "Generating Sequences With RNNs," 2014)

GPT-3 showed some ability to continue relatively complex patterns. But only with GPT-3.5—able to handle multi-stage dialogue and follow elaborate instructions—it became clear that something unexpected happens when a language model surpasses a certain parameter scale and is pretrained on a sufficiently large corpus.

Scale is fundamental to building a capable LLM. Let's look at the core features that make LLMs "large" and how these features contribute to their capabilities.

### 5.1.1. Large Parameter Count

One of the most striking features of LLMs is the sheer number of parameters they contain. While our decoder model has around 8 million parameters, state-of-the-art LLMs can reach hundreds of billions or even trillions of parameters.

In a transformer model, the number of parameters is largely determined by the embedding dimensionality (`emb_dim`) and the number of decoder blocks (`num_blocks`). As these values increase, the parameter count grows quadratically with embedding dimensionality in the self-attention and MLP layers, and linearly with the number of decoder blocks. For example, doubling the embedding dimensionality roughly quadruples the number of parameters in the attention and MLP components of each decoder block.

**Open-weight models** are models with publicly accessible trained parameters. These can be downloaded and used for tasks like text generation or finetuned for specific applications. However, while the weights are open, the model's license governs its permitted uses, including whether commercial use is allowed. Licenses like Apache 2.0 and MIT permit unrestricted commercial use, but you should always review the license to confirm your intended use aligns with the creators' terms.

The table below shows key features of several open-weight LLMs compared to our tiny model:

	<code>num_blocks</code>	<code>emb_dim</code>	<code>num_heads</code>	<code>vocab_size</code>
Our model	2	128	8	32,011
Llama 3.1 8B	32	4,096	32	128,000
Gemma 2 9B	42	3,584	16	256,128
Gemma 2 27B	46	4,608	32	256,128
Llama 3.1 70B	80	8,192	64	128,000
Llama 3.1 405B	126	16,384	128	128,000

By convention, the number before “B” in the name of an open-weight model indicates its total number of parameters in billions.

If you were to store each parameter of a 70B model as a 32-bit float number, it would require about 280GB of RAM—more storage than the Apollo 11 guidance computer had by a factor of over 30 million times.

This massive number of parameters allows LLMs to learn and represent a vast amount of information about grammar, semantics, world knowledge, and exhibit reasoning capabilities.

### 5.1.2. Large Context Size

Another crucial aspect of LLMs is their ability to process and maintain much larger contexts than earlier models. While our decoder model used a context of only 30 tokens, modern LLMs can handle contexts of thousands—and sometimes even millions—of tokens.

GPT-3’s 2,048-token context could accommodate roughly 4 pages of text. In contrast, Llama 3.1’s 128,000-token context is large enough to fit the entire text of *“Harry Potter and the Sorcerer’s Stone”* with room to spare.

The key challenge with processing long texts in transformer models lies in the self-attention mechanism’s computational complexity. For a sequence of length  $n$ , self-attention requires computing attention scores between every pair of tokens, resulting in quadratic  $O(n^2)$  time and space complexity. This means that doubling the input length quadruples both the memory requirements and computational cost. This quadratic scaling becomes particularly problematic for long documents—for instance, a 10,000-token input would require computing and storing 100 million attention scores for each attention layer.

The increased context size is made possible through architectural improvements and optimizations in attention computation. Techniques like **grouped-query attention** and **FlashAttention** (which are beyond the scope of this book) enable efficient memory management, allowing LLMs to handle much larger contexts without excessive computational costs.

LLMs typically undergo pretraining on shorter contexts around 4K-8K tokens, as the attention mechanism’s quadratic complexity makes training on long sequences computationally intensive. Additionally, most training data naturally consists of shorter sequences.

Long-context capabilities emerge through **long-context pretraining**, a specialized stage following initial training. This process involves:

1. **Incremental training for longer contexts:** The model progressively adapts from 4K-8K tokens to 128K-256K tokens through staged increases. Each stage continues until the model demonstrates both preserved performance on short-context tasks and success on specific evaluations like “needle in a haystack” tests.

**A needle in a haystack** test evaluates a model's ability to identify and utilize relevant information buried within a very long context, typically by placing a crucial piece of information early in the sequence and asking a question that requires retrieving that specific detail from among thousands of tokens of unrelated text.

2. **Efficient scaling for self-attention:** To handle the computational demands of self-attention's quadratic scaling with sequence length, the approach implements **context parallelism**. This method splits input sequences into manageable chunks and uses an all-gather mechanism for memory-efficient processing.

**All-gather** is a collective communication operation in distributed computing where each GPU shares its local data with all other GPUs, aggregating the data so that every GPU ends up with a complete, concatenated dataset. This operation is particularly useful in transformer models, where different GPUs process parts of a sequence in parallel but require access to the full context for tasks like attention computation.

While Llama 2 and 3 models maintain a constant  $\theta$  of 500,000 throughout pretraining, Qwen 2.5 models employ a dynamic approach during long-context pretraining, gradually increasing  $\theta$  from 10,000 to 1,000,000.

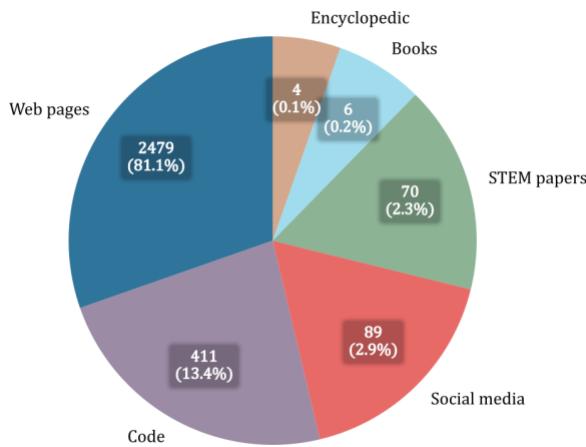
A higher base frequency in RoPE reduces rotation angles, decreasing attention score decay across long distances and enabling better long-range information aggregation. The increased base frequency also provides finer granularity in representing relative distances between positional embeddings, allowing the model to better differentiate positions in long contexts.

### 5.1.3. Large Training Dataset

The third factor behind LLMs' capabilities is the size of the corpus used for training. While our decoder was trained on a small corpus of news sentences with about 25 million tokens, modern LLMs use datasets with trillions of tokens. These datasets often include:

- 1) books and literature from different genres and eras,
- 2) web pages and online articles on diverse topics,
- 3) academic papers and scientific studies,
- 4) social media posts and discussions, and
- 5) code repositories and technical documents.

The diversity and scale of these datasets allow LLMs to learn a broad vocabulary, understand multiple languages, acquire knowledge on a wide array of topics—from history and science to current events and pop culture—adapt to various writing styles and formats, and acquire basic reasoning and problem-solving skills.



The illustration depicts the composition of LLM training datasets, using the open **Dolma** dataset as an example. Segments represent different document types, with sizes scaled logarithmically to prevent web pages—the largest category—from overwhelming the visualization. Each segment shows both token count (in billions) and percentage of the corpus. While Dolma’s 3 trillion tokens are substantial, they fall short of more recent datasets like Qwen 2.5’s 18 trillion tokens, a number likely to grow in future iterations.

It would take approximately 51,000 years for a human to read the entire Dolma dataset, reading 8 hours every day at 250 words per minute.

Since neural language models train on such vast corpora, they typically process the data just once. This **single-epoch training** approach prevents **overfitting** while reducing computational demands. Processing these enormous datasets multiple times would be extremely time-consuming and may not yield significant additional benefits.

#### 5.1.4. Large Amount of Compute

If you tried to process 3 trillion tokens of the Dolma dataset on a single modern GPU, it would take over 100 years—which helps explain why major language models require massive computing clusters. Training an LLM demands significant computing power, often measured in **FLOPs (floating-point operations)** or GPU-hours. For context, while training our decoder model might take a few hours on a single GPU, modern LLMs can require thousands of GPUs running for months.

The computational demands grow with three main factors:

- 1) the number of parameters in the model,
- 2) the size of the training corpus, and
- 3) the context length used during training.

For example, training the Llama 3.1 series of models consumed approximately 40 million GPU-hours—equivalent to running a single GPU continuously for almost 4600 years. Llama 3.1’s training process uses an advanced system called **4D parallelism**, which integrates four different parallel processing methods to efficiently distribute the model across thousands of GPUs.

The first dimension is **tensor parallelism**, which splits individual weight matrices across devices. This includes attention weight matrices ( $\mathbf{W}^Q$ ,  $\mathbf{W}^K$ ,  $\mathbf{W}^V$ , and  $\mathbf{W}^O$ ) and MLP matrices ( $\mathbf{W}_1$ ,  $\mathbf{W}_2$ ). **Pipeline parallelism** forms the second dimension, distributing different transformer layers across GPUs, with each GPU handling specific layers of the model. The third dimension, **context parallelism**, divides input sequences into chunks for efficient processing of long sequences. Finally, **data parallelism** allows multiple GPUs to process different batches simultaneously, synchronizing after each training step.

Each of these four parallelism dimensions could merit its own chapter, and thus a full exploration of them lies beyond this book’s scope.

Training large language models can cost tens to hundreds of millions of dollars. These expenses include hardware, electricity, cooling, and engineering expertise. Such costs limit the development of state-of-the-art LLMs to large tech companies and well-funded research labs. However, open-weight models have lowered the barrier, enabling smaller organizations to leverage existing models through methods like supervised finetuning and prompt engineering.

## 5.2. Supervised Finetuning

During pretraining, the model learns most of its capabilities. However, since it is trained only to predict the next word, its default behavior is to continue the input. For instance, if you input “Explain how machine learning works,” the pretrained model might respond with something like “and also name three most popular algorithms.” This is not what users would expect. The model’s ability to follow instructions, answer questions, and hold conversations is developed through a process called **supervised finetuning**.

Let’s compare the behavior of a pretrained model and the same model finetuned to follow instructions and answer questions.

We’ll use two models: `google/gemma-2-2b`, pretrained for next-token prediction, and `google/gemma-2-2b-it`, a finetuned version for instruction following.

Models on the **Hugging Face Hub** follow this naming convention: “creator/model” with no spaces. The “model” part typically includes information about the model’s version, number of parameters, and whether it was finetuned for conversation or instruction-following. In the name `google/gemma-2-2b-it`, we see that the creator is Google, the model has version 2, 2 billion parameters, and it was finetuned to follow instructions (with “it” standing for “instruction-tuned”).

The prompt we’ll use is the following:

Continue the list: apple, orange, tomato.

This is the output of the pretrained-only google/gemma-2-2b given the above prompt:

The list of fruits and vegetables that are good for you is long. But there are some that are better than others.

The best fruits and vegetables are those that are high in fiber, low in sugar, and high in vitamins and minerals.

The best fruits and vegetables are those that are high in fiber, low in sugar, and high in vitamins and minerals.

...

The output isn't complete—the model keeps repeating the same sentence endlessly. As you can see, the output is quite similar to what we observed with our decoder model. While google/gemma-2-2b, being larger, produces more coherent sentence structures, the text still fails to align with the context, which clearly requests a list of fruits.

Now, let's apply the finetuned google/gemma-2-2b-it to the same input. The output is:

\* \*\*Fruits\*\*

Here are some more fruits to continue the list:

\* \*\*Banana\*\*  
\* \*\*Strawberry\*\*  
\* \*\*Grape\*\*  
\* \*\*Pineapple\*\*  
\* \*\*Mango\*\*  
\* \*\*Blueberry\*\*  
\* \*\*Peach\*\*  
\* \*\*Watermelon\*\*  
\* \*\*Lemon\*\*  
\* \*\*Lime\*\*

Let me know if you'd like to add more!

As you can see, the model with the same number of parameters now follows the instruction. This change is achieved through supervised finetuning.

**Supervised finetuning**, or simply **finetuning**, modifies a pretrained model's parameters to specialize it for specific tasks. The goal isn't to train the model to answer every question or follow every instruction. Instead, finetuning "unlocks" the knowledge and skills the model already learned during pretraining. Without finetuning, this knowledge remains "hidden" and is used mainly for predicting the next token, not for problem-solving.

While the model continues to operate as a next-token predictor during finetuning, the key difference lies in the training data. Instead of predicting tokens from general text, the model now learns from examples of high-quality conversations, instruction-following, and problem-solving. This specialized training helps the model learn to apply its preexisting knowledge in more useful ways, revealing its learned facts and patterns in response to prompts rather than simply continuing arbitrary text.

### 5.3. Finetuning a Pretrained Model

As discussed, training an LLM from scratch is a complex and expensive undertaking that requires significant computational resources, vast amounts of high-quality training data, as well as deep expertise in machine learning research and engineering.

The good news is that you don't need to train these models from scratch to solve business problems. Open-weight models often come with permissive licenses that allow you to use or finetune them for specific tasks without buying them or paying royalties. Models with up to 8 billion parameters can be finetuned even in a Colab notebook (in the paid version that supports more powerful GPUs), though the process will be slow, and a single GPU may lack sufficient memory to store both the model and long prompts.

To speed up finetuning and process longer contexts, organizations often use servers with multiple high-end GPUs running in parallel. Each GPU has substantial VRAM (video random access memory), which stores models and data during computation. By distributing the model's weights across the GPUs' combined memory, finetuning becomes significantly faster than relying on a single GPU. This approach is called **model parallelism**.

PyTorch supports model parallelism with methods like **Fully Sharded Data Parallel (FSDP)**. FSDP enables efficient distribution of model parameters across GPUs by **sharding** the model—splitting it into smaller parts. This way, each GPU processes only a portion of the model.

These servers can be quite expensive for smaller organizations or individuals. The cost per GPU-hour typically ranges from \$2 to \$4, and finetuning a model can take anywhere from hours to weeks, depending on the model size and the dataset used.

Commercial LLM service providers offer a more cost-effective option for finetuning large models. They charge based on the number of tokens in the training data and use various techniques to lower costs. Though these methods aren't covered in this book, you can find an up-to-date list of LLM finetuning services with pay-per-token pricing on the book's wiki.

Let's finetune a pretrained LLM to predict an emotion. Our dataset has the following structure:

```
{"text": "i slammed the door and screamed in rage", "label": "anger"}  
 {"text": "i danced and laughed under the bright sun", "label": "joy"}  
 {"text": "tears rolled down my face in silence today", "label": "sadness"}  
 ...
```

It's a JSONL file, where each row is a labeled example as a JSON object. The `text` key contains a text expressing one of six emotions; the `label` key is the corresponding emotion. The label can be one of six values: sadness, joy, love, anger, fear, and surprise. Thus, we have a document classification problem with six classes.

We'll finetune GPT-2, a pretrained model licensed under the MIT license, which permits unrestricted commercial use. GPT-2 is one of the smallest LLMs, often referred to as an SLM (small language model). Despite its smaller size, it performs well for certain tasks and can be finetuned in a Colab notebook.

Before training a complex model, it's wise to establish baseline performance. A **baseline** is a simple, easy-to-implement solution that sets the minimum acceptable performance level. Without it, we can't determine if a complex model's performance justifies its added complexity.

As a baseline, we'll use **logistic regression** combined with **bag of words**. This combination is well-known for producing effective models in document classification tasks. We will use the open-source machine learning library **scikit-learn** that simplifies training and evaluation of traditional machine learning models.

### 5.3.1. Baseline Emotion Classifier

First, we install scikit-learn:

```
$ pip3 install scikit-learn
```

Now, let's load the data and prepare it for machine learning:<sup>10</sup>

```
random.seed(42) ❶
data_url = "https://www.thelmbbook.com/data/emotions"
X_train_text, y_train, X_test_text, y_test = download_and_split_data(
    data_url, test_ratio=0.1
) ❷
```

The function `download_and_split_data` (defined in the [thelmbbook.com/nb/5.1](#) notebook) downloads a compressed dataset from a specified URL, extracts the training examples, and splits the dataset into **training** and **test** partitions. The `test_ratio` parameter in line ❷ specifies the fraction of the dataset to reserve for testing. Setting a seed in ❶ ensures that the random shuffle in line ❷ produces the same result on every execution for reproducibility.

With the data loaded and split into training and test sets, we transform it into a bag-of-words:

```
from sklearn.feature_extraction.text import CountVectorizer
```

---

<sup>10</sup> We will load the data from the book's website to ensure it remains accessible. The dataset's original source is <https://huggingface.co/datasets/dair-ai/emotion>. It was first used in Saravia et al., "CARER: Contextualized Affect Representations for Emotion Recognition," Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, 2018.

```
vectorizer = CountVectorizer(max_features=10_000, binary=True)
X_train = vectorizer.fit_transform(X_train_text)
X_test = vectorizer.transform(X_test_text)
```

The `fit_transform` method of `CountVectorizer` transforms the training data into a bag-of-words representation. The `max_features` parameter controls vocabulary size, while `binary` specifies whether features are binary; if `False`, non-zero values indicate word counts. The `transform` method in the next line converts the test data into a bag-of-words using the vocabulary created by `fit_transform` on the training data. This avoids **data leakage**, where test data unintentionally affects training. A key principle in machine learning is that test data must never influence model training.

The logistic regression implementation in scikit-learn accepts labels as strings, so there is no need to convert them to numbers. The library handles the conversion automatically.

Now, let's train a logistic regression model:

```
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score

model = LogisticRegression(random_state=42, max_iter=1000)
model.fit(X_train, y_train) # Model is trained here

y_train_pred = model.predict(X_train)
y_test_pred = model.predict(X_test)

train_accuracy = accuracy_score(y_train, y_train_pred)
test_accuracy = accuracy_score(y_test, y_test_pred)

print(f"Training accuracy: {train_accuracy * 100:.2f}%")
print(f"Test accuracy: {test_accuracy * 100:.2f}%")
```

Output:

```
Training accuracy: 0.9854
Test accuracy: 0.8855
```

The `LogisticRegression` object is first created. Its `fit` method, called next, trains the model<sup>11</sup> on the training data. Afterward, the model predicts outcomes for both the training and test sets, and the accuracy for each is calculated.

The `random_state` parameter in `LogisticRegression` sets the seed for the random number generator. The `max_iter` parameter limits the solver to a maximum of 1000 iterations to achieve **convergence**. A **solver** is the algorithm used to optimize the model's parameters. This ensures the solver has sufficient iterations to find the optimal solution.

---

<sup>11</sup> In reality, scikit-learn trains a model slightly different from classical logistic regression; it uses softmax with cross-entropy loss instead of the sigmoid function. This approach enables multiclass classification.

The **accuracy** metric calculates the proportion of correct predictions out of all predictions:

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

As you can see, the model **overfits**: it performs almost perfectly on the training data but significantly worse on the test data. To address this, we can adjust the **hyperparameters** of our algorithm. Let's try incorporating bigrams and increase the vocabulary size to 20,000:

```
vectorizer = CountVectorizer(max_features=20_000, ngram_range=(1, 2))
```

This adjustment leads to slight improvement on the test set, but it still falls short compared to the training set performance:

```
Training accuracy: 0.9962
Test accuracy: 0.8910
```

Now that we see a simple approach achieves a test accuracy of 0.8910, any more complex solution must outperform this baseline. If it performs worse, we will know that our implementation probably contains an error.

Let's finetune GPT-2 to generate emotion labels as text. This approach is easy to implement since no additional classification output layer is needed. Instead, the model is trained to output labels as regular words, which may span multiple tokens.

### 5.3.2. Emotion Generation

First, we get the data, model, and tokenizer:

```
from transformers import AutoTokenizer, AutoModelForCausalLM

set_seed(42)
data_url = "https://www.thelmbbook.com/data/emotions"
model_name = "openai-community/gpt2"
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
tokenizer = AutoTokenizer.from_pretrained(model_name) ❶
tokenizer.pad_token = tokenizer.eos_token ❷
model = AutoModelForCausalLM.from_pretrained(model_name).to(device) ❸
num_epochs, batch_size, learning_rate = get_hyperparameters()
train_loader, test_loader = download_and_prepare_data(
    data_url, tokenizer, batch_size
)
```

The `AutoModelForCausalLM` class from the `transformers` library, used in line ❸, automatically loads a pretrained **autoregressive language model**. Line ❶ loads the pretrained tokenizer. The tokenizer used in GPT-2 does not include a padding token. Therefore, in line ❷, we set the padding token by reusing the end-of-sequence token.

Now, we set up the training loop:

```

for epoch in range(num_epochs):
    for input_ids, attention_mask, labels in train_loader:
        input_ids = input_ids.to(device)
        attention_mask = attention_mask.to(device) ❶
        labels = labels.to(device)
        outputs = model(
            input_ids=input_ids,
            labels=labels,
            attention_mask=attention_mask
        )
        outputs.loss.backward()
        optimizer.step()
        optimizer.zero_grad()

```

The `attention_mask` in line ❶ is a binary tensor showing which tokens in the input are actual data and which are padding. It has 1s for real tokens and 0s for padding tokens. This mask is different from the **causal mask**, which blocks positions from attending to future tokens.

Let's illustrate `input_ids`, `labels`, and `attention_mask` for a batch of two simple examples:

Text	Emotion
I feel very happy	joy
So sad today	sadness

We convert these examples into text completion tasks by adding a task definition and solution:

Table 5.1: Text completion template.

Task	Solution
Predict emotion: I feel very happy\nEmotion:	joy
Predict emotion: So sad today\nEmotion:	sadness

In the table above, “\n” denotes a new line character, while “\nEmotion:” marks the boundary between the task description and the solution. This format, while optional, helps the model use its pretrained understanding of text. The sole new ability to be learned during finetuning is generating one of six outputs: sadness, joy, love, anger, fear, or surprise—no other outputs.

LLMs gained emotion classification skills during pretraining partly because of the widespread use of emojis online. Emojis acted as labels for the text around them.

Assuming a simple tokenizer that splits strings by spaces and assigns unique IDs to each token, here's a hypothetical token-to-ID mapping:

Token	ID	Token	ID
Predict	1	So	8

<b>Token</b>	<b>ID</b>	<b>Token</b>	<b>ID</b>
emotion:	2	sad	9
I	3	today	10
feel	4	joy	11
very	5	sadness	12
happy	6	[EOS]	0
\nEmotion:	7	[PAD]	-1

The special [EOS] token indicates the end of generation, while [PAD] serves as a padding token. The following examples show how texts are converted to token IDs:

<b>Text</b>	<b>Token IDs</b>
Predict emotion: I feel very happy\nEmotion: joy	[1, 2, 3, 4, 5, 6, 7] [11]
Predict emotion: So sad today\nEmotion: sadness	[1, 2, 8, 9, 10, 7] [12]

We then concatenate the input tokens with the completion tokens and append the [EOS] token so the model learns to stop generating once the emotion label generation is completed. The `input_ids` tensor directly contains these concatenated token IDs. The `labels` tensor is made by replacing all input text tokens with -100 (a special masking value), while keeping the actual token IDs for the completion and [EOS] tokens. This ensures the model only computes loss on predicting the completion tokens, not on reproducing the input text.

The value -100 is a special token ID in PyTorch (and similar frameworks) used to exclude specific positions during loss computation. When finetuning language models, this ensures the model concentrates on predicting tokens for the desired output (the “solution”) rather than the tokens in the input (the “task”).

Here’s the resulting table:

<b>Text</b>	<b>input_ids</b>	<b>labels</b>
Predict emotion: I feel very happy\nEmotion: joy	[1, 2, 3, 4, 5, 6, 7, 11, 0]	[-100, -100, -100, -100, -100, -100, -100, 11, 0]
Predict emotion: So sad today\nEmotion: sadness	[1, 2, 8, 9, 10, 7, 12, 0]	[-100, -100, -100, -100, -100, -100, -100, 12, 0]

To form a batch, all sequences must have the same length. The longest sequence has 9 tokens (from the first example), so we pad the shorter sequences to match that length. Here’s the final table showing how the `input_ids`, `labels`, and `attention_mask` are adjusted after padding:

input_ids	labels	attention_mask
[1, 2, 3, 4, 5, 6, 7, 11, 0]	[-100, -100, -100, -100, -100, -100, -100, 11, 0]	[1, 1, 1, 1, 1, 1, 1, 1, 1]
[1, 2, 8, 9, 10, 7, 12, 0, -1]	[-100, -100, -100, -100, -100, -100, -100, 12, 0, -100]	[1, 1, 1, 1, 1, 1, 1, 1, 0]

In `input_ids`, all sequences are padded to a length of 9 tokens. The second example is padded with the `[PAD]` token (ID  $-1$ ). In the `attention_mask`, real tokens are marked as 1, while padding tokens are marked as 0.

This padded batch is now ready for the model to handle.

After finetuning the model with `num_epochs = 2`, `batch_size = 16`, and `learning_rate = 0.00005`, it achieves a test accuracy of 0.9415. This is more than 5 percentage points higher than the baseline result of 0.8910 obtained with logistic regression.

When finetuning, a smaller learning rate is usually applied to prevent drastic changes to the model's pretrained weights. This approach helps preserve the general knowledge learned during pretraining while adapting to the new task. A learning rate of  $0.00005$  ( $5 \times 10^{-5}$ ) is a common choice, as it has been empirically found to work well for many transformer-based models and tasks, though optimal values may vary depending on the task and model.

The full code for supervised finetuning of an LLM is available in the [thelmbook.com/nb/5.2](https://huggingface.co/course/notebooks/5.2_finetuning.ipynb) notebook. You can adapt this code for any text generation task by updating the data files (while keeping the same JSON format) and adjusting Task and Solution in Table 5.1 with text relevant to the specific business problem.

Let's see how this code can be adapted for finetuning for a general instruction-following task.

### 5.3.3. Finetuning to Follow Instructions

While similar to the emotion generation task, let's quickly review the specifics of finetuning a large language model to follow arbitrary instructions.

When finetuning a language model for instruction-following, the first step is choosing a **prompting format** or **prompting style**. For emotion generation, we used this format:

```
Predict emotion: {text}
Emotion: {emotion}
```

This format allows the LLM to see where the Task part ends (“\nEmotion:”) and the Solution starts. When we finetune for a general-purpose instruction following, we cannot use “\nEmotion:” as a separator. We need a more general format. Since first open-weight models were introduced, many prompting formats were used by various people and organizations. Below, there are only some of them, named after famous LLMs using these formats:

Vicuna:

```
USER: {instruction}
ASSISTANT: {solution}
```

Alpaca:

```
### Instruction:
{instruction}

### Response:
{solution}
```

**ChatML (chat markup language)** is a prompting format used in many popular finetuned LLMs. It provides a standardized way to encode chat messages, including the role of the speaker and the content of the message.

The format uses two tags: <|im\_start|> to indicate the start of a message and <|im\_end|> to mark its end. A basic ChatML message structure looks like this:

```
<|im_start|>{role}
{message}
<|im_end|>
```

The **message** is either an instruction (question) or a solution (answer). The **role** is usually one of the following: **system**, **user**, and **assistant**. For example:

```
<|im_start|>system
You are a helpful assistant.
<|im_end|>
<|im_start|>user
What is the capital of France?
<|im_end|>
<|im_start|>assistant
The capital of France is Paris.
<|im_end|>
```

The **user** role is the person who asks questions or gives instructions. The **assistant** role is the chat LM providing responses. The **system** role specifies instructions or context for the model's behavior. The **system** message, known as the **system prompt**, can include private details about the user, like their name, age, or other information useful for the LLM-based application.

The prompting format has little impact on the quality of a finetuned model itself. However, when working with a model finetuned by someone else, you need to know the format used during finetuning. Using the wrong format could affect the quality of the model's outputs.

After transforming the training data into the chosen prompting format, the training process uses the same code as the emotion generation model. You can find the complete code for instruction finetuning an LLM in the [thelmbbook.com/nb/5.3](http://thelmbbook.com/nb/5.3) notebook.

The dataset I used contains about 500 examples, generated using a state-of-the-art LLM. While 500 examples may not be sufficient for achieving high-quality instruction following, there's no universally accepted method for creating an ideal instruction finetuning dataset. Online datasets range from hundreds of thousands to millions of examples. Yet some experiments indicate that with a well-chosen set of examples spanning diverse instruction types, even 1,000 examples can enable strong instruction-following capabilities in a sufficiently large pretrained language model.

A consensus among the practitioners is that the quality, not quantity, of examples is crucial for achieving state-of-the-art results in instruction finetuning.

The training examples can be found in this file:

```
data_url = "https://www.thelmbbook.com/data/instruct"
```

It has the following structure:

```
...
{"instruction": "Translate 'Good night' into Spanish.", "solution": "Buenas noches"}
 {"instruction": "Name three primary colors.", "solution": "Red, blue, yellow"}
 ...
...
```

The instructions and examples used during finetuning fundamentally shape a model's behavior. When training data contains many polite or cautious responses, the finetuned model tends to produce similar answers. Through finetuning, a model can even be trained to consistently generate falsehoods. Users working with third-party finetuned LLMs should remain mindful of potential biases introduced in the process.

To understand the impact of instruction finetuning, let's first see how a pretrained model handles instructions without any special training. Let's first use a pretrained GPT-2:

```
from transformers import AutoTokenizer, AutoModelForCausalLM
import torch

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

tokenizer = AutoTokenizer.from_pretrained("openai-community/gpt2")
tokenizer.pad_token = tokenizer.eos_token

model = AutoModelForCausalLM.from_pretrained("openai-community/gpt2").to(device)

instruction = "Who is the President of the United States?"
inputs = tokenizer(instruction, return_tensors="pt").to(device)

outputs = model.generate(
```

```

    input_ids=inputs["input_ids"],
    attention_mask=inputs["attention_mask"],
    max_new_tokens=32,
    pad_token_id=tokenizer.pad_token_id
)

generated_text = tokenizer.decode(outputs[0], skip_special_tokens=True)
print(generated_text)

```

Output:

```
Who is the President of the United States?
```

```
The President of the United States is the President of the United States.
```

```
The President of the United States is the President of the United States.
```

Again, google/gemma-2-2b, the model exhibits sentence repetition. Now, let's look at the output after finetuning on our instruction dataset. The inference code for an instruction-finetuned model must follow the prompting format used during finetuning. The `build_prompt` method applies the ChatML prompting format to our instruction:

```

def build_prompt(instruction, solution = None):
    wrapped_solution = ""
    if solution:
        wrapped_solution = f"\n{n{solution}}\n<|im_end|>"
    return f""""<|im_start|>system
You are a helpful assistant.
<|im_end|>
<|im_start|>user
{instruction}
<|im_end|>
<|im_start|>assistant"""+wrapped_solution

```

The same `build_prompt` function is used for both training and testing. During training, it takes both `instruction` and `solution` as input. During testing, it only receives `instruction`.

Now, let's define the function that generates text:

```

def generate_text(model, tokenizer, prompt, max_new_tokens=100):
    input_ids = tokenizer(prompt, return_tensors="pt").to(model.device)

    end_tokens = tokenizer.encode("<|im_end|>", add_special_tokens=False) ①

    stopping = [EndTokenStoppingCriteria(end_tokens, model.device)] ②

    output_ids = model.generate(
        input_ids=input_ids["input_ids"],
        attention_mask=input_ids["attention_mask"],

```

```

        max_new_tokens=max_new_tokens,
        pad_token_id=tokenizer.pad_token_id,
        stopping_criteria=stopping
    )[0]

    generated_ids = output_ids[input_ids["input_ids"].shape[1]:] ③
    generated_text = tokenizer.decode(generated_ids).strip()
    return generated_text

```

Line ① encodes the `<|im_end|>` tag into token IDs which will be used to indicate the end of generation. Line ② sets up a stopping criterion using the `EndTokenStoppingCriteria` class (defined below), ensuring the generation halts when `end_tokens` appear. Line ③ slices the generated tokens to remove the input prompt, leaving only the newly generated text.

The `EndTokenStoppingCriteria` class defines the signal to stop generating tokens:

```

class EndTokenStoppingCriteria(StoppingCriteria):
    def __init__(self, end_tokens, device):
        self.end_tokens = torch.tensor(end_tokens).to(device) ①

    def __call__(self, input_ids, scores):
        do_stop = []
        for sequence in input_ids: ②
            if len(sequence) >= len(self.end_tokens):
                last_tokens = sequence[-len(self.end_tokens):] ③
                do_stop.append(torch.all(last_tokens == self.end_tokens)) ④
            else:
                do_stop.append(False)
        return torch.tensor(do_stop, device=input_ids.device)

```

In the constructor:

- Line ① converts the `end_tokens` list into a PyTorch tensor and moves it to the specified device. This ensures the tensor is on the same device as the model.

In the `__call__` method, line ② loops through the generated sequences in the batch. For each:

- Line ③ takes the last `len(end_tokens)` tokens and stores them in `last_tokens`.
- Line ④ checks if `last_tokens` matches `end_tokens`. If they do, `True` is added to the `do_stop` list, which tracks whether to stop generation for each sequence in the batch.

This is how we call the inference for a new instruction:

```

input_text = "Who is the President of the United States?"
prompt = build_prompt(input_text)
generated_text = generate_text(model, tokenizer, prompt)
print(generated_text.replace("<|im_end|>", "").strip())

```

Output:

George W. Bush

Since GPT-2 is a relatively small language model and wasn't finetuned on recent facts, this confusion about presidents isn't surprising. What matters here is that the finetuned model now interprets the instruction as a question and responds accordingly.

## 5.4. Sampling From Language Models

When generating text with a language model, we need to turn the output logits into tokens. **Greedy decoding**, where we select the highest probability token at each step, works well for tasks like math or factual questions that require precision. However, many applications benefit from some randomness. For example, brainstorming story ideas works better when the model generates diverse outputs. Debugging code can improve with alternative suggestions if the first one fails. Even for summarization or translation, sampling helps explore equally valid phrasing options when the model is uncertain.

To address this, we *sample* from the probability distribution instead of always choosing the most likely token. Different techniques allow us to control how much randomness to introduce.

Let's explore some of these techniques.

### 5.4.1. Basic Sampling with Temperature

The simplest approach converts logits to probabilities using the **softmax** function with a **temperature** parameter  $T$ :

$$\Pr(j) = \frac{\exp(o^{(j)}/T)}{\sum_{k=1}^V \exp(o^{(k)}/T)}$$

where  $o^{(j)}$  represents the logit for token  $j$ ,  $\Pr(j)$  gives its resulting probability, and  $V$  denotes the vocabulary size. The temperature  $T$  determines the sharpness of the probability distribution:

- At  $T = 1$ , we obtain standard softmax probabilities.
- As  $T \rightarrow 0$ , the distribution focuses on the highest probability tokens.
- As  $T \rightarrow \infty$ , the distribution approaches uniformity.

For example, if we have logits  $[4, 2, 0]^\top$  for tokens "cat", "dog", and "bird" (assuming only three words in the vocabulary), here's how different temperatures affect the probabilities:

$T$	Probabilities	Comment
0.5	$[0.98, 0.02, 0.00]^\top$	More focused on "cat"
1.0	$[0.87, 0.12, 0.02]^\top$	Standard softmax
2.0	$[0.67, 0.24, 0.09]^\top$	More evenly distributed

Temperature controls the trade-off between creativity and determinism. Lower temperatures (such as 0.1–0.3) typically produce more focused and conservative outputs, which may be

preferred for tasks requiring precision like factual responses, code generation, or math. Moderate temperatures (often around 0.7–0.8) can balance creativity and coherence, making them useful starting points for conversation or content writing. Higher temperatures (commonly 1.5–2.0) introduce more randomness, which can help with tasks like brainstorming or story generation, though outputs may become less coherent as temperature increases. While these ranges serve as general guidelines, the optimal temperature varies significantly by model, task, and specific implementation. Most applications avoid extreme temperatures (very close to 0 or above 2), but the best value should be determined through experimentation.

Given the vocabulary and probabilities, this Python function returns the sampled token:

```
import numpy as np

def sample_token(probabilities, vocabulary):
    if len(probabilities) != len(vocabulary): ❶
        raise ValueError("Mismatch between the two inputs' sizes.")

    if not np.isclose(sum(probabilities), 1.0, rtol=1e-5): ❷
        raise ValueError("Probabilities must sum to 1.")

    return np.random.choice(vocabulary, p=probabilities) ❸
```

The function performs two checks before sampling. Line ❶ ensures there is one probability for each token in the vocabulary. Line ❷ confirms the probabilities sum to 1, allowing for a small tolerance due to floating-point precision. Once these validations pass, line ❸ handles the sampling. It selects a token from the vocabulary based on the probabilities, so a token with a 0.7 probability is chosen roughly 70% of the time when the function is run repeatedly.

#### 5.4.2. Top- $k$ Sampling

While temperature helps control randomness, it allows sampling from the entire vocabulary, including very unlikely tokens that the model assigns extremely low probabilities to. **Top- $k$  sampling** addresses this by limiting the sampling pool to the  $k$  most likely tokens.

The process works as follows:

1. Sort tokens by probability,
2. Keep only the top  $k$  tokens,
3. Renormalize their probabilities to sum to 1,
4. Sample from this reduced distribution.

We can update `sample_token` to support both temperature and top- $k$  sampling:

```
def sample_token(logits, vocabulary, temperature=0.7, top_k=50):
    if len(logits) != len(vocabulary):
        raise ValueError("Mismatch between logits and vocabulary sizes.")
    if temperature <= 0:
        raise ValueError("Temperature must be positive.")
```

```

if top_k < 1:
    raise ValueError("top_k must be at least 1.")
if top_k > len(logits):
    raise ValueError("top_k must be at most len(logits).")

logits = logits / temperature ●

cutoff = np.sort(logits)[-top_k] ●
logits[logits < cutoff] = float("-inf") ●

probabilities = np.exp(logits - np.max(logits)) ●
probabilities /= probabilities.sum() ●

return np.random.choice(vocabulary, p=probabilities)

```

The function begins by validating inputs: ensuring logits match the vocabulary size, temperature is positive, top-k is at least 1, and top-k does not exceed the vocabulary size. Line ❶ scales the logits by the temperature. Line ❷ determines the top-k cutoff by sorting the logits and selecting the  $k^{\text{th}}$  largest value. Line ❸ discards less likely tokens by setting logits below the cutoff to negative infinity. Line ❹ converts the remaining logits into probabilities using a numerically stable softmax. Line ❺ ensures the probabilities sum to 1.

Subtracting `np.max(logits)` before exponentiating avoids numerical overflow. Large logits can produce excessively large exponentials. Shifting the largest logit to 0 keeps values stable while preserving their relative proportions.

The value of  $k$  varies by task. Lower values (5–10) prioritize the most probable tokens, enhancing accuracy and consistency—beneficial for factual responses and structured tasks. Mid-range values (20–50) strike a balance between variation and coherence, serving as effective defaults for general writing and dialogue. Higher values (100–500) incorporate less likely tokens to enable greater diversity, particularly useful for creative tasks like brainstorming. Though these ranges offer practical guidance, optimal  $k$  selection depends heavily on the model, vocabulary size, and specific application. Extremely low values (under 5) can be overly restrictive, while very high values (above 500) seldom improve output quality. Like temperature, finding the ideal value requires experimentation.

#### 5.4.3. Nucleus (Top-p) Sampling

**Nucleus sampling**, or **top-p sampling**, takes a different approach to token selection. Instead of using a fixed number of tokens, it selects the smallest group of tokens whose cumulative probability exceeds a threshold  $p$ .

Here's how it works for  $p = 0.9$ :

1. Rank tokens by probability,
2. Add tokens to the subset until their cumulative probability surpasses 0.9,

3. Renormalize the probabilities of this subset,
4. Sample from the adjusted distribution.

This method adapts to the context. It might select just a few tokens for highly focused distributions or many tokens when the model is less certain.

In practice, these three methods are often used together in the following sequence:

1. **Temperature scaling** (e.g.,  $T = 0.7$ ) adjusts the randomness by sharpening or softening the probabilities of tokens.
2. **Top-k filtering** (e.g.,  $k = 50$ ) limits the sampling pool to the  $k$  most probable tokens, ensuring computational efficiency and preventing extremely low-probability tokens from being considered.
3. **Top-p filtering** (e.g.,  $p = 0.9$ ) further refines the sampling pool by selecting the smallest set of tokens whose cumulative probability meets the threshold  $p$ .

#### 5.4.4. Penalties

In addition to temperature and filtering methods, modern language models use several penalty parameters to control the diversity and quality of generated text. These parameters help prevent common issues like repetition of words or phrases, overuse of common tokens, and getting stuck in generation loops.

The **frequency penalty** adjusts token probabilities based on how often they've appeared in the generated text so far. When a token appears multiple times, its probability is reduced proportionally to its appearance count. The penalty is applied by subtracting a scaled version of the token's count from its logits before the softmax:

$$o^{(j)} \leftarrow o^{(j)} - \alpha \cdot \text{count}(j),$$

where  $\alpha$  is the frequency penalty parameter. Higher values (0.8-1.0) decrease the model's likelihood to repeat the same line verbatim or getting stuck in a loop.

The **presence penalty** modifies token probabilities based on whether they appear anywhere in the generated text, regardless of count:

$$o^{(j)} \leftarrow \begin{cases} o^{(j)} - \gamma, & \text{if token } j \text{ is in generated text,} \\ o^{(j)}, & \text{otherwise} \end{cases}$$

Here,  $\gamma$  is the presence penalty parameter. Higher values of  $\gamma$  (0.7-1.0) increase the model's likelihood to talk about new topics.

The optimal values depend on the specific task. For creative writing, higher penalties encourage novelty. For technical documentation, lower penalties maintain precision and consistency.

The complete implementation of `sample_token` that combines temperature, top- $k$ , top- $p$ , and the two penalties can be found in the [thelmbbook.com/nb/5.4](https://thelmbbook.com/nb/5.4) notebook.

## 5.5. Low-Rank Adaptation (LoRA)

Finetuning LLMs through adjustment of their billions of parameters requires extensive computational resources and memory, creating barriers for those with limited infrastructure.

**LoRA (low-rank adaptation)** offers a solution by updating fewer parameters during finetuning. The method introduces small matrices to capture necessary adjustments rather than modifying the entire model. This approach achieves comparable performance while training only a fraction of the parameters.

### 5.5.1. The Core Idea

In the Transformer, most parameters are found in the weight matrices of **self-attention** and **position-wise MLP** layers. Rather than modifying the large weight matrices directly, LoRA introduces two smaller matrices for each. During finetuning, these smaller matrices are trained to capture the required adjustments, while the original weight matrices stay “frozen.”

Consider a  $d \times k$  weight matrix  $\mathbf{W}_0$  in a pretrained model. Instead of updating  $\mathbf{W}_0$  directly during finetuning, we modify the process like this:

1. **Freeze the original weights:** The matrix  $\mathbf{W}_0$  remains unchanged during finetuning.
2. **Add two small matrices:** Introduce an  $d \times r$  matrix  $\mathbf{A}$  and an  $r \times k$  matrix  $\mathbf{B}$ , where  $r$ —referred to as the **rank**—is an integer much smaller than both  $d$  and  $k$  (e.g.,  $r = 8$ ).
3. **Adjust the weights:** Compute the adapted weight matrix  $\mathbf{W}$  during finetuning as:

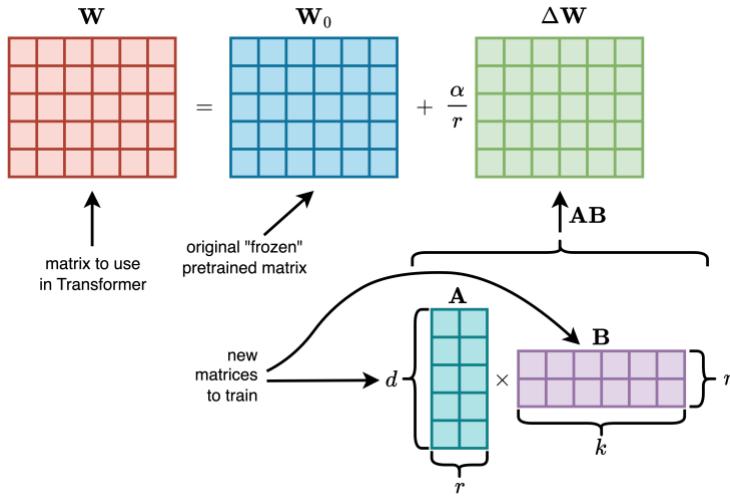
$$\mathbf{W} = \mathbf{W}_0 + \frac{\alpha}{r} \Delta\mathbf{W} = \mathbf{W}_0 + \frac{\alpha}{r} \mathbf{AB}$$

Here,  $\Delta\mathbf{W} = \mathbf{AB}$  represents the adjustment to  $\mathbf{W}_0$ , scaled by the **scaling factor**  $\frac{\alpha}{r}$ .

The matrices  $\mathbf{A}$  and  $\mathbf{B}$  are called a **LoRA adapter**. Their product,  $\Delta\mathbf{W}$ , acts as an update matrix that adjusts the original weights  $\mathbf{W}_0$  to enhance performance on a new task. Since  $\mathbf{A}$  and  $\mathbf{B}$  are much smaller than  $\mathbf{W}_0$ , this method significantly reduces the number of trainable parameters.

For example, if  $\mathbf{W}_0$  has dimensions  $1024 \times 1024$ , it would contain over a million parameters to finetune directly (1,048,576 parameters). With LoRA, we introduce  $\mathbf{A}$  with dimensions  $1024 \times 8$  (8,192 parameters) and  $\mathbf{B}$  with dimensions  $8 \times 1024$  (8,192 parameters). This setup requires only  $8,192 + 8,192 = 16,384$  parameters to be trained.

The adapted weight matrix  $\mathbf{W}$  is used in the layers of the finetuned transformer, replacing the original matrix  $\mathbf{W}_0$  to alter the token embeddings as they pass through the transformer blocks. The creation of  $\mathbf{W}$  is shown below:



The scaling factor  $\frac{\alpha}{r}$  controls the size of the weight updates introduced by LoRA during finetuning. Both  $r$  and  $\alpha$  are hyperparameters, with  $\alpha$  typically set as a multiple of  $r$ . For example, if  $r = 8$ ,  $\alpha$  might be 16, resulting in a scaling factor of 2. The optimal values for  $r$  and  $\alpha$  are found experimentally by assessing the finetuned LLM's performance on the test set.

LoRA is applied to the weight matrices in the self-attention layers—specifically the query, key, and value weight matrices  $\mathbf{W}^Q$ ,  $\mathbf{W}^K$ ,  $\mathbf{W}^V$ , and the **projection matrix**  $\mathbf{W}^O$ . It can also be applied to the weight matrices  $\mathbf{W}_1$  and  $\mathbf{W}_2$  in the position-wise MLP layers.

Finetuning LLMs with LoRA is faster than a **full model finetune** and uses less memory for gradients, enabling the finetuning of very large models on limited hardware.

### 5.5.2. Parameter-Efficient Finetuning (PEFT)

The Hugging Face **Parameter-Efficient Finetuning (PEFT)** library provides a simple way to implement LoRA in transformer models. Let's install it first:

```
pip3 install peft
```

We can modify our previous code by incorporating the PEFT library to apply LoRA:

```
from peft import get_peft_model, LoraConfig, TaskType

peft_config = LoraConfig(
    task_type=TaskType.CAUSAL_LM, # Specify the task type
    inference_mode=False, # Set to False for training
    r=8, # Set the rank
    lora_alpha=16 # LoRA alpha
)
```

```
model = get_peft_model(model, peft_config)
```

The `LoraConfig` object defines the parameters for LoRA finetuning:

- `task_type` specifies the task, which in this case is **causal language modeling**,
- `r` is the LoRA adapter rank,
- `lora_alpha` is the scaling factor  $\alpha$ .

We use `get_peft_model` to wrap the original model and integrate LoRA adapters. You might wonder how it decides which matrices to augment with LoRA adapters. PEFT is built to recognize standard LLM architectures. So, when finetuning models like Llama, Gemma, Mistral, or Qwen, it automatically applies LoRA to the relevant layers. For custom transformers—like our decoder from Chapter 4—you can include an extra parameter, `target_modules`, to specify which matrices should use LoRA:

```
peft_config = LoraConfig(
    #same as above
    target_modules=["W_Q", "W_K", "W_V", "W_O"]
)
```

Next, we set up the optimizer as usual:

```
optimizer = torch.optim.AdamW(model.parameters(), lr=learning_rate)
```

In PyTorch, the `requires_grad` attribute controls whether a tensor tracks operations for automatic differentiation. When `requires_grad=True`, PyTorch keeps track of all operations on the tensor, enabling gradient computation during the backward pass. To freeze a model parameter (preventing updates during training), set its `requires_grad` to `False`:

```
import torch.nn as nn

model = nn.Linear(2, 1) # Linear Layer: y = Wx + b

print(model.weight.requires_grad)
print(model.bias.requires_grad)

model.bias.requires_grad = False
print(model.bias.requires_grad)
```

Output:

```
True
True
False
```

The PEFT library ensures that only the LoRA adapter parameters have `requires_grad=True`, keeping all other model parameters frozen.

After wrapping the model with `get_peft_model`, the training loop stays the same. For instance, finetuning GPT-2 on an emotion generation task using LoRA with `r=16` and `lora_alpha=32` achieves a test accuracy of 0.9420. This is marginally better than the 0.9415 from full finetuning. Generally, LoRA tends to perform slightly worse than full finetuning. However, the outcome depends on the choice of hyperparameters, dataset size, base model, and task.

The full code for GPT-2 finetuning with LoRA is available in the [thelmbbook.com/nb/5.5](https://thelmbbook.com/nb/5.5) notebook. You can customize it for your own tasks by modifying the dataset and LoRA settings.

## 5.6. LLM as a Classifier

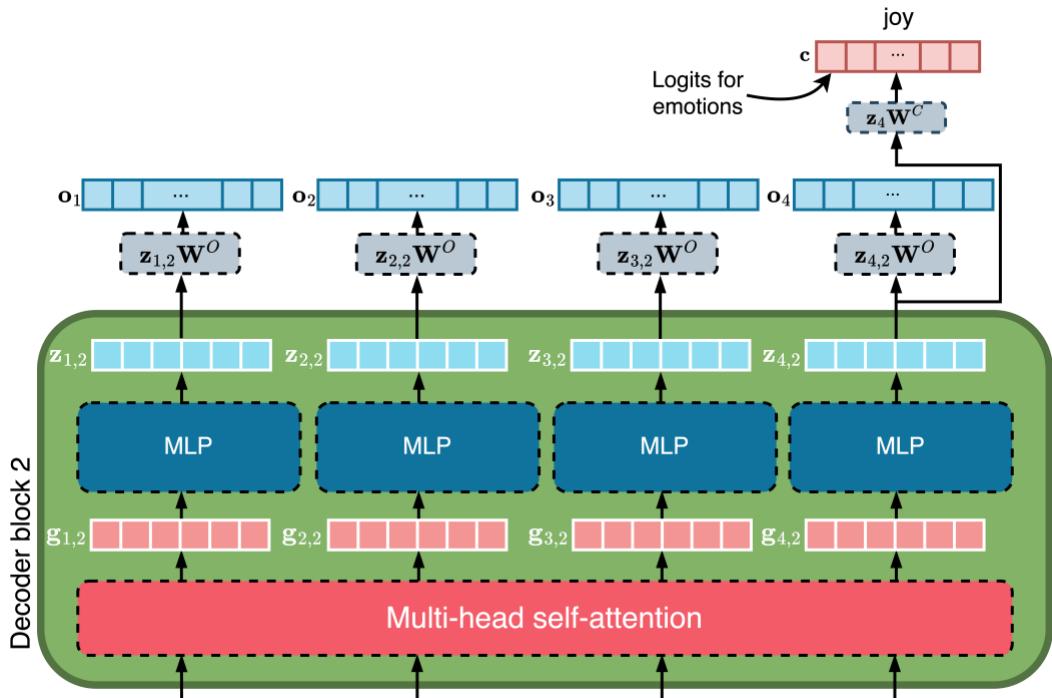
When finetuning GPT-2 for emotion prediction, we didn't turn it into a classifier. Instead, it generated the class name as text. While this method works, it's not always optimal for classification tasks. A different approach is to train the model to produce logits for each emotion class.

We can attach a **classification head** to a pretrained LLM. This is a fully connected layer with a softmax activation. The softmax maps logits to probabilities, and the class with the highest probability becomes the predicted label.

In `transformers`, there is a class designed to make this easier. Instead of loading the model with `AutoModelForCausalLM`, we use `AutoModelForSequenceClassification`:

```
from transformers import AutoModelForSequenceClassification  
  
model = AutoModelForSequenceClassification.from_pretrained(model_path, num_labels=6)
```

For pretrained autoregressive language models, the class maps the embedding of the final (right-most) non-padding token from the last decoder block to a vector with dimensions matching the number of classes (6 in this case). The structure of this modification is as follows:



As you can see, once the final decoder block processes the input (the second block in our example), the output embedding  $\mathbf{z}_{4,2}$  of the last token is passed through the classification head's weight matrix,  $\mathbf{W}^C$ . This projection converts the embedding into logits, one per class.

The parameter tensor  $\mathbf{W}^C$  is initialized with random values and trained on the labeled emotions dataset. Training relies on **cross-entropy** to measure the loss between the predicted probability distribution and the **one-hot encoded** true class label. This error is backpropagated, updating the weights in both the classification head and the rest of the model.

After finetuning with `num_epochs = 8`, `batch_size = 16`, and `learning_rate = 0.00005`, the model reaches a test accuracy of 0.9460. This is slightly better than the 0.9415 accuracy from finetuning the unmodified model to generate class labels as text. The improvement might be more noticeable with a different base model or dataset.

The code for finetuning GPT-2 as an emotion classifier is available on the wiki in the [thelmbbook.com/nb/5.6](https://thelmbbook.com/nb/5.6) notebook. It can be easily adapted for any text classification task by replacing the data in the file while keeping the same JSON format.

## 5.7. Prompt Engineering

**Chat language models**, or **chat LMs**, are language models finetuned on dialogue examples. This finetuning resembles instruction finetuning but uses multi-turn conversation inputs, such as those in the ChatML format, with the targets being the assistant's responses.

Despite its simplicity, the conversational interface allows solving various practical problems. This section explores best practices for using chat LMs to address such problems known as **prompt engineering** techniques.

### 5.7.1. Features of a Good Prompt

To get the best results from a chat LM, you need a well-crafted prompt. The key components of a strong prompt include:

1. **Situation:** Describe why you're asking for help.
2. **Role:** Define the expert persona the model should emulate.
3. **Task:** Give clear, specific instructions about what the model must do.
4. **Output format:** Explain how you expect the response to be structured, such as bullet points, JSON, or code.
5. **Constraints:** Mention any limitations, preferences, or requirements.
6. **Quality criteria:** Define what makes a response satisfactory.
7. **Examples:** Provide few-shot examples of inputs with expected outputs.
8. **Call to action:** Restate the task simply and ask the model to perform it.

Putting input-output examples in the prompt is called **few-shot prompting** or **in-context learning**. These examples include both positive cases showing desired outputs and negative ones demonstrating incorrect responses. Adding explanations that connect errors to specific constraints helps model to understand why they are wrong.

Here's an example of a prompt that includes some of the above elements:

Situation: I'm creating a system to analyze insurance claims. It processes adjuster reports to extract key details for display in a SaaS platform.

Your role: Act as a seasoned insurance claims analyst familiar with industry-standard classifications.

Task: Identify the type of incident, the primary cause, and the significant damages described in the report.

Output format: Return a JSON object with this structure:

```
{  
    "type": "string",      // Incident type  
    "cause": "string",     // Primary cause  
    "damage": ["string"]   // Major damages  
}
```

```
<examples>  
    <example>  
        <input>  
            Observed two-vehicle accident at an intersection. Insured's car was  
            hit after the other driver ran a red light. Witnesses confirm. The vehicle  
            has severe front-end damage, airbags deployed, and was towed from the scene
```

```
.  
    </input>  
    <output>  
    {  
        "type": "collision",  
        "cause": "failure to stop at signal",  
        "damage": ["front-end damage", "airbag deployment"]  
    }  
    </output>  
  </example>  
  <example>  
    ...  
  </example>  
</examples>
```

Call to action: Extract the details from this report:

"Arrived at the scene of a fire at a residential building. Extensive damage to the kitchen and smoke damage throughout. Fire caused by unattended cooking. Neighbors evacuated; no injuries reported."

Section names such as "Situation," "Your role," or "Task" are optional.

I used XML tags for few-shot examples because they clearly define example boundaries and are familiar to LLMs from pretraining on structured data. Furthermore, chat LM models are often finetuned using conversational examples with XML structures. Using XML isn't mandatory though, but could be helpful.

The attention mechanism in LLMs has limitations. It might concentrate on certain parts of a prompt while overlooking others. A good prompt strikes a balance between detail and brevity. Excessive detail can overwhelm the model, while insufficient detail risks leaving gaps that the model may fill with incorrect assumptions.

### 5.7.2. Followup Actions

The first solution from a model is often imperfect. User analysis and follow-up are key to getting the most out of a chat LM. Common follow-up actions include:

1. Asking the LLM whether its solution contains errors or can be simplified without breaking the constraints.
2. Copying the solution and starting a new conversation from scratch with the same LLM. In this new conversation, the user can ask the model to validate the solution as if it were "provided by an expert," without revealing it was generated by the same model.
3. Using a different LLM to review or enhance the solution.
4. For code outputs, running the code in the execution environment, analyzing the results, and giving feedback to the model. If code fails, the full error message and stack traceback can be shared with the model.

When working with the same chat LM for follow-ups, especially in tasks like coding or handling complex structured outputs, it's generally a good idea to start fresh after three-five exchanges. This recommendation comes from two key observations:

1. Chat LMs are typically finetuned using examples of short conversations. Creating long, high-quality conversations for finetuning is both difficult and costly, so the training data often lacks examples of long interactions focused on problem solving. As a result, the model performs better with shorter exchanges.
2. Long contexts can cause errors to accumulate. In the self-attention mechanism, the softmax is applied over many positions to compute weights for combining value vectors. As the context length increases, inaccuracies build up, and the model's "focus" may shift to irrelevant details or earlier mistakes.

When starting fresh, it's important to update the initial prompt with key details from earlier followups. This helps the model avoid repeating previous mistakes. By consolidating the relevant information into a clear, concise starting point, you ensure the model has the context it needs without relying on the history of the prior conversation.

### 5.7.3. Code Generation

One valuable use of chat LMs is generating code. The user describes the desired code, and the model tries to generate it. As we know, modern LLMs are pretrained on vast collections of open-source code across many programming languages. This pretraining allows them to learn syntax and many standard or widely-used libraries. Seeing the same algorithms implemented in different languages also enables LLMs to form shared internal representations, making them generally indifferent to the programming language when reading or creating code.

Moreover, much of this code includes comments and annotations, which help the model understand the code's purpose—what it is designed to achieve. Sources like StackOverflow and similar forums add further value by providing examples of problems paired with their solutions. The exposure to such data gave LLMs an ability to respond with relevant code. Supervised finetuning improved their skill in interpreting user requests and turning them into code.

As a result, LLMs can generate code in nearly any language. For high-quality results, users must specify in detail what the code should do. For example, providing a detailed docstring like this:

```
Write Python code that implements a method with the following specification  
s:
```

```
def find_target_sum(numbers: list[int], target: int) -> tuple:  
    """Find pairs of indices in a list whose values sum to a target.
```

Args:

    numbers: List of integers to search through. Can be empty.  
    target: Integer sum to find.

Returns:

    Tuple of two distinct indices whose values sum to target,  
    or None if no solution exists.

**Examples:**

```
>>> find_target_sum([2, 7, 11, 15], 9)
(0, 1)
>>> find_target_sum([3, 3], 6)
(0, 1)
>>> find_target_sum([1], 5)
None
>>> find_target_sum([], 0)
None
```

**Requirements:**

- Time complexity:  $O(n)$
- Space complexity:  $O(n)$
- Each index can only be used once
- If multiple solutions exist, return any valid solution
- All numbers and target can be any valid integer
- Return None if no solution exists

....

Providing a highly detailed docstring can sometimes feel as time-consuming as coding the function itself. A less detailed description might seem more practical, but this increases the likelihood of the generated code not fully meeting user needs. In such cases, users can review the output and refine their instructions with additional requests or constraints.

Language models can generate functions, classes, or even entire applications. However, the chance of success decreases as the level of abstraction increases. If the problem resembles model's training data, the model performs well with minimal input. However, for novel or unique business problems, detailed instructions are crucial for good results.

If you use a brief prompt to save time, ask the model to pose clarifying questions. You can also request it to describe the code it plans to generate first. This allows you to adjust or add details to the instructions before code is created.

## 5.8. Hallucinations

A major challenge with LLMs is their tendency to produce content that seems plausible but is factually incorrect. These inaccuracies, called **hallucinations**, create problems for using LLMs in production systems where reliability and accuracy are required.

### 5.8.1. Reasons for Hallucinations

Hallucinations in LLMs are by design. As we know, these models are optimized to predict the next token that fits the context, not to ensure factual accuracy. During pretraining, they learn to generate coherent text by following language patterns. However, no training dataset can cover every fact. When the model faces knowledge gaps, it tries to fill them with plausible content based on patterns it recognizes. This often results in fabricated details. Here's an example of hallucinated information from a widely used chat LM:

Explain the principle of blockchain quantum neural network.

A Blockchain Quantum Neural Network (BQNN) combines three cutting-edge technologies: **blockchain**, **quantum computing**, and **neural networks**. It aims to leverage the strengths of each to create secure, efficient, and intelligent systems. Here's a breakdown of the principle:

As you can imagine, “Blockchain Quantum Neural Network (BQNN)” is not a real concept. The LLM’s two-page explanation, including detailed descriptions of how it works, is entirely fabricated.

Low quality of training data also contributes to hallucinations. During pretraining on large volumes of internet text, models are exposed to both accurate and inaccurate information. They learn these inaccuracies but lack the ability to differentiate between truth and falsehood.

Finally, LLMs generate text one token at a time. This approach means that errors in earlier tokens can cascade, leading to increasingly incoherent outputs.

### 5.8.2. Preventing Hallucinations

Hallucinations cannot be completely avoided, but they can be minimized. A practical way to reduce hallucinations is by grounding the model’s responses in verified information. This is done by including relevant context directly in the prompt. For instance, rather than posing an open-ended question, we can provide specific documents or data for the model to reference and instruct the model to only answer based on the provided documents.

This method, called **retrieval-augmented generation (RAG)**, anchors the model’s output to verifiable facts. The model still generates text but does so—most of the time—with the limits of the provided context, which significantly reduces hallucinations.

Here’s how RAG works: when a user submits a query, the system searches a knowledge base, like a document repository or database, for relevant information. It combines keyword matching with embedding-based search, where the user input is transformed into an embedding vector, and documents with similar embeddings are retrieved, often using **cosine similarity**. To prevent overwhelming the model with lengthy documents, the documents are split into smaller chunks before embedding.

The retrieved content is added to the prompt alongside the user’s question. This approach merges the strengths of traditional information retrieval with the language generation capabilities of LLMs. For example, if a user asks about a company’s latest quarterly results, the RAG system would first retrieve the most recent financial reports and use them to produce the response, avoiding reliance on potentially outdated training data.

Another way to reduce hallucinations is by finetuning the model on reliable, domain-specific knowledge using unlabeled documents. For instance, a question-answering system for law firms

could be finetuned on legal documents, case law, and statutes to improve accuracy within the legal domain. This approach is often referred to as **domain-specific pretraining**.

For critical applications, implementing a multi-step verification workflow can provide additional protection against hallucinations. This might involve using multiple models with different architectures or training data to cross-validate responses and having domain experts review generated content before it's used in production.

However, it's important to recognize that hallucinations cannot be completely eliminated with current LLM technology. While we can implement various safeguards and detection mechanisms, the most robust approach is to design systems that account for this limitation.

For instance, in a customer service application, an LLM could draft responses, but human review would be necessary before sending messages containing specific product details or policy information. Similarly, in a code generation system, the model might generate code, but automated tests and human review should always occur before deployment.

The potential for hallucinations was notably demonstrated when Air Canada's customer service chatbot provided incorrect information about bereavement travel rates to a passenger. The chatbot falsely claimed that customers could book full-price tickets and later apply for reduced fares, contradicting the airline's actual policy. When the passenger tried to claim the fare reduction, Air Canada's denial led to a small claims court case, resulting in an \$812 CAD (near \$565 USD) compensation order. This case highlights the tangible business consequences of AI inaccuracies, including financial losses, customer frustration, and reputational damage.

Success with LLMs lies in recognizing that hallucinations are an inherent limitation of the technology. However, this issue can be managed through thoughtful system design, safeguards, and a clear understanding of when and where these models should be applied.

## 5.9. LLMs, Copyright, and Ethics

The widespread deployment of LLMs has introduced novel challenges in copyright law, particularly regarding training data usage and the status of AI-generated content. These issues affect both the companies developing LLMs and the businesses building applications with them.

### 5.9.1. Training Data

The first major copyright consideration involves training data. LLMs are trained on large text datasets that include copyrighted material such as books, articles, and software code. While some claim that this might qualify as fair use,<sup>12</sup> this has not been tested in court. The issue is further complicated by the models' capacity to output protected content. This legal uncertainty has already

<sup>12</sup> Fair use is a U.S. legal doctrine. Other regions handle copyright exceptions differently. The EU relies on "fair dealing" and specific statutory exceptions, Japan has distinct copyright limitations, and other countries apply unique rules for permitted uses. This variation complicates global LLM deployment, as training data allowed under U.S. fair use might violate copyright laws elsewhere.

sparked high-profile lawsuits from authors and publishers against AI companies, posing risks for businesses using LLM applications.

Meta's decision to withhold its multimodal Llama model from the European Union in July 2024 exemplifies the growing tension between AI development and regulatory compliance. Citing concerns over the region's "unpredictable" regulatory environment, particularly regarding the use of copyrighted and personal data for training, Meta joined other tech giants like Apple in limiting AI deployments in European markets. This restriction highlights the challenges companies face in balancing innovation with regional regulations.

When selecting models for commercial use, companies should review the training documentation and license terms. Models trained primarily on public domain or properly licensed materials involve lower legal risks. However, the massive datasets required for effective LLMs make it nearly impossible to avoid copyrighted material entirely. Businesses need to understand these risks and factor them into their development strategies.

Beyond legal issues, training LLMs on copyrighted material raises ethical concerns. Even when legally permissible, using copyrighted works without consent may appear exploitative, especially if the model outputs compete with the creators' work. Transparency about training data sources and proactive engagement with creators can help address these concerns. Ethical practices should also involve compensating creators whose contributions significantly improve the model, fostering a more equitable system.

### 5.9.2. Generated Content

The copyright status of content generated by LLMs presents challenges that traditional copyright law cannot easily resolve. Copyright law is built around the assumption of human authorship, leaving it unclear whether AI-generated works qualify for protection or who the rightful owner might be. Another issue is that LLMs can sometimes reproduce portions of their training data verbatim, including copyrighted material. This ability to generate exact reproductions—beyond learning abstract patterns—raises serious legal questions.

Some businesses address these challenges by treating LLMs as tools to assist humans rather than as independent creators. For instance, a marketing team might use an LLM to draft text, which human writers then edit and finalize. This ensures clearer copyright ownership while still benefiting from AI's efficiency. Similarly, software developers rely on LLMs to generate code snippets, which they review and incorporate into larger systems. This practice has become widespread—at Google, over 25% of all code written in 2024 was generated by LLMs, then reviewed and refined by developers.

To minimize copyright risks in LLM applications, companies often implement technical safeguards.

One method involves comparing model outputs against a database of copyrighted materials to detect verbatim copies. This process identifies direct matches or significant similarities using techniques like hash-based comparisons or advanced NLP algorithms. For example, a company may maintain a repository of copyrighted texts and employ similarity detection methods—such as cosine similarity or edit distance—to flag outputs that surpass a defined similarity threshold.

However, these methods are not foolproof. Paraphrased content can make the output formally distinct while remaining substantively similar, which automated systems may fail to detect. To handle this, businesses often supplement these tools with human review to ensure compliance.

### 5.9.3. Open-Weight Models

The copyright status of model weights introduces legal questions distinct from those related to training data or generated outputs. Model weights, which encode the patterns learned during training, might be seen as derivative works of the training data. This raises the issue of whether sharing weights constitutes an indirect redistribution of the original copyrighted materials, even in their transformed state. Some argue that weights are an abstract transformation and represent entirely new intellectual property. Others believe that if weights enable the reproduction of training data fragments, they inherently incorporate copyrighted content and should be subject to similar legal protections.

This debate carries serious implications for open-source AI development. If model weights are classified as derivative works, sharing and distributing models trained on copyrighted data could become legally restricted, even if the training process qualifies as fair use. As a result, some organizations have shifted to training models solely on public domain or explicitly licensed content. However, this strategy often limits the effectiveness of the models, as the smaller, restricted datasets typically lead to reduced performance.

As laws around LLMs evolve, businesses must stay flexible. They may need to adjust workflows as courts define legal boundaries or revise policies as AI-specific legislation appears. Consulting intellectual property lawyers with AI expertise can help manage these shifts.

# Chapter 6. Further Reading

You've learned the core concepts of language models throughout this book. There are many advanced topics to explore on your own, and this final chapter provides pointers for further study. I've chosen topics that represent important current developments in the field, from architectural innovations to security considerations.

## 6.1. Mixture of Experts

**Mixture of experts (MoE)** is an architectural pattern designed to increase model capacity without a proportional rise in cost. Instead of a single position-wise MLP processing all tokens in a decoder block, MoE uses multiple specialized sub-networks called **experts**. A **router network** (or **gate network**) decides which tokens are processed by which experts.

The core idea is activating only a subset of experts for each token. This **sparse** activation reduces active computations while enabling larger overall parameter counts. **Sparse MoE layers** replace traditional MLP layers, using techniques like **top-k routing** and **load balancing** to efficiently assign tokens to experts.

This concept gained attention with the **Switch Transformer** and has been applied in models such as **Mixtral 8x7B**. For example, Mixtral achieves 47B total parameters but only activates about 13B during inference.

## 6.2. Model Merging

**Model merging** combines multiple pretrained models to make use of their complementary strengths. Techniques include **model soups**, **SLERP** (spherical interpolation that maintains parameter norms), and **task vector algorithms** such as **TIES-Merging** and **DARE**.

These methods generally rely on some architectural similarity or compatibility between models. The **passthrough** method stands out by concatenating layers from different LLMs. This approach can create models with unconventional parameter counts (e.g., 13B by merging two 7B models). Such models are often called **frankenmerges**.

**mergekit** is a popular open-source tool for merging and combining language models that implements many of these techniques. It provides a flexible configuration system for experimenting with different merging strategies and architectures.

## 6.3. Model Compression

**Model compression** addresses deploying LLMs in resource-limited environments by reducing size and computation needs without greatly sacrificing performance. Neural networks are often **over-parameterized**, containing redundant units that can be optimized. Key methods include **post-training quantization**, which lowers parameter precision (e.g., 32-bit floats to 8-bit integers), **quantization-aware training**, which trains models at lower precision, such as **QLoRA** (quantized low-rank adaptation), **unstructured pruning**, removing individual weights by importance, **structured pruning**, removing components like layers or attention heads, and **knowledge distillation**, where a smaller "student" model learns from a larger "teacher" model.

## 6.4. Preference-Based Alignment

**Preference-based alignment** methods help align LLMs with user values and intent so they produce helpful and safe outputs. A widely used approach is **reinforcement learning from human feedback (RLHF)**, where humans rank model responses, a **reward model** is trained on these rankings, and then the LLM is finetuned to optimize for higher reward. Another approach is **constitutional AI (CAI)**, which uses a set of guiding principles or a “constitution” that the model refers to when producing its output; the model can **self-critique** and revise its responses based on these principles. Both strategies address the problem that LLMs, when trained on vast internet text, may generate harmful or **misaligned** responses, but they differ in how they incorporate human oversight and explicit guidelines.

## 6.5. Advanced Reasoning

Advanced reasoning techniques enable large language models to handle complex tasks by (1) training them to generate an explicit **chain of thought (CoT)** for step-by-step reasoning and (2) equipping them with **function calling** capabilities to invoke external APIs or tools, thereby addressing limitations of simple prompt-response patterns. Chain-of-thought reasoning can significantly improve performance on tasks such as multi-step mathematics and logical inference, while function calling allows offloading specialized computations to external frameworks.

Additionally, **tree of thought (ToT)** extends CoT by exploring multiple reasoning paths in a tree-like structure. **Self-consistency** further refines reasoning by aggregating multiple CoT outputs for the most consistent answer. **ReAct (reasoning+act)** integrates reasoning with action-taking, allowing models to interact with environments dynamically. **Program-aided language models (PAL)** leverage interpreters (e.g., Python) to execute code for precise calculations.

## 6.6. Language Model Security

**Jailbreak attacks** and **prompt injection** are major security vulnerabilities in LLMs. Jailbreaks bypass the model’s safety controls by crafting specific inputs that trick the model into producing restricted content, often using techniques like roleplaying as a different character or setting up hypothetical scenarios. For example, an attacker might prompt the model to act as a pirate to obtain instructions on illegal activities. In contrast, prompt injection attacks manipulate how LLM applications combine **system prompts** with user input, allowing attackers to alter the application’s behavior. For instance, an attacker could insert commands that make the application execute unauthorized actions. While jailbreaks primarily risk exposing harmful or restricted content, prompt injection presents more severe security implications for applications with privileged access, such as those that read emails or execute system commands.

## 6.7. Vision Language Model

**Vision language models (VLMs)** integrate an LLM with a **vision encoder** to handle both text and images. Unlike traditional models that process modalities in isolation, VLMs excel at **multimodal reasoning**, enabling them to perform a variety of vision tasks by following natural language instructions without task-specific retraining. The architecture includes three main components: a **CLIP**-based (contrastive language-image pretraining) **vision encoder** trained on millions of image-

text pairs to understand visual content, a **cross-attention** mechanism that allows the VLM to integrate and reason about visual and textual information, and the language model itself that generates and interprets text. VLMs are developed through multiple training stages, starting with pretraining to align the visual and language components, followed by supervised finetuning to improve their ability to understand and respond to user prompts.

## 6.8. Preventing Overfitting

Techniques for preventing **overfitting** are essential for achieving model **generalization**, ensuring that models perform well not just on training data but also on new, unseen examples. The primary defense against overfitting is **regularization**, which includes methods like **L1** and **L2**. These techniques add specific penalty terms—such as the sum of absolute or squared weights—to the loss function, limiting the size of model parameters and encouraging simpler models. **Dropout** is a powerful regularization method for neural networks that randomly deactivates a portion of units during each training step. This forces the network to develop multiple independent pathways, preventing units from relying too heavily on specific features. **Early stopping** offers a practical approach by tracking the model’s performance on a validation set and stopping training when the validation accuracy stops improving or begins to decline, thereby avoiding the model from memorizing random noise in the training data.

A **validation set** is similar to the **test set** in that it is used to evaluate the model’s performance on unseen data; however, the key difference is that the validation set is used during the training process to tune hyperparameters and make decisions such as early stopping, while the test set is reserved for final evaluation to measure the model’s performance after training is complete.

## 6.9. Concluding Remarks

You’ve come a long way in understanding language models, from the basic building blocks of machine learning to the inner workings of transformers and the practical aspects of working with large language models. You now have a solid technical foundation that lets you not only understand how these models work but also implement and adapt them for your own purposes.

The field of language models continues to evolve, with new architectures, training methods, and applications emerging. You’re now equipped to read research papers, understand technical discussions, and evaluate new developments critically. Whether you plan to develop models, finetune existing ones, or build systems that use them, you have the core concepts needed to proceed confidently.

I encourage you to stay curious and hands-on—implement the concepts you’ve learned, experiment with different approaches, and keep up with the latest developments. Consider starting with some of the advanced topics covered in this chapter, but remember that the fundamentals you’ve learned here will serve as your compass in navigating future innovations.