

# LLM Basics I

- An LLM learns a model that uses the current pre-context to predict the next token.
- We give the following example:

| Pre-context             |   | The next token |
|-------------------------|---|----------------|
| I                       | → | am             |
| I am                    | → | a              |
| I am a                  | → | machine        |
| I am a machine          | → | learning       |
| I am a machine learning | → | researcher     |

# LLM Basics II

- This setting is called **autoregressive prediction** in machine learning.
- Assume we have  $D$  documents. For document  $i$ , let us assume

$\mathbf{x}_{i,1:j}$ : contexts from 1st to  $j$ th words

$\mathbf{y}_{ij}$ :  $(j + 1)$ st word.

- Then we have many  
(target, instance)  
pairs for training.

- For example, we have

$\mathbf{x}_i$ : “I am a machine learning researcher” and

# LLM Basics III

$\mathbf{x}_{i,1:1}$  |  
 $\mathbf{x}_{i,1:2}$  | am  
 $\vdots$

- From the basic concept of supervised learning, we can solve an optimization problem

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^D \sum_j \xi(f(\boldsymbol{\theta}; \mathbf{x}_{i,1:j}), \mathbf{y}_{ij})$$

to obtain the model, where

- $\boldsymbol{\theta}$  is the model parameters,

# LLM Basics IV

- $f(\theta; \mathbf{x}_{i,1:j})$  is a function that gives an approximation of the target  $\mathbf{y}_{ij}$ ,
- $\xi(\cdot)$  is the loss function.
- However, training a supervised learning model requires that

$$\mathbf{y}_{ij} \text{ and } \mathbf{x}_{i,1:j}$$

are **vectors**, but they are now word sequences.

# LLM Basics V

- Having vectors is important because, for example, we can then define a simple loss function such as

$$\|\mathbf{y}_{ij} - f(\boldsymbol{\theta}; \mathbf{x}_{i,1:j})\|^2$$

to ensure that  $\boldsymbol{\theta}$  leads to an output as close to  $\mathbf{y}_{ij}$  as possible.

- From what we learned in supervised learning, if instances and targets are **vectors**, and we agree that the following procedures can be implemented:
  - network architectures (i.e., our  $f$  function),
  - stochastic gradient methods,

# LLM Basics VI

- automatic differentiation,  
then we “believe” that an LLM can be implemented.
- We will explain how we can convert  $(\mathbf{y}_{ij}, \mathbf{x}_{i,1:j})$  to vectors.
- Subsequently, we abuse the notation a bit so that  $\mathbf{x}_{i,1:j}$  is a word sequence as well as the converted vector.

# Tokenization I

- We split each document to several tokens and map each token to an integer ID.
- An example is as follows:

|                |                                    |           |          |                |                 |                   |      |
|----------------|------------------------------------|-----------|----------|----------------|-----------------|-------------------|------|
| Document       | I am a machine learning researcher |           |          |                |                 |                   |      |
|                |                                    |           |          |                | ↓               |                   |      |
| Tokenized text | <u>I</u>                           | <u>am</u> | <u>a</u> | <u>machine</u> | <u>learning</u> | <u>researcher</u> |      |
|                |                                    |           |          |                | ↓               |                   |      |
| Token IDs      |                                    | 25        | 7        | 100            | 2               | 10                | 1000 |

# Tokenization II

- We let “Vocabulary” be the huge dictionary of all words (tokens) considered.
- Thus,  $|\text{Vocabulary}|$  is its size.
- For easy understanding, we consider a word as a token.
- In practice, we often break a word to several tokens.



# Tokenization III

- If each word is a token, for every word in the sentence we can check the corresponding ID in the Vocabulary dictionary:

the **25th** word in our dictionary is “I”,

the **7th** word in our dictionary is “am”,

the **100th** word in our dictionary is “a”,

the **2nd** word in our dictionary is “machine”,

the **10th** word in our dictionary is “learning”,

the **1000th** word in our dictionary is “researcher”.

# Position Information I

- We can let

$$\mathbf{x}_{i,1:j} \in \{0, 1\}^{|\text{Vocabulary}|} \quad (1)$$

be an indicator vector to reflect the occurrence of its words.

- For example, if

$$\mathbf{x}_{i,1:2} \text{ is "I am"} \quad (2)$$

# Position Information II

then we have

$$\mathbf{x}_{i,1:2} = \underbrace{[0, \dots, 0, 1, 0, \dots, 0, 1, 0, \dots]}_{\text{7th}}^{\text{25th}}]^\top.$$

- However, from the indicator vector and the Vocabulary, we can extract only the set of words in  $\mathbf{x}_{i,1:j}$  instead of the word sequences.
- That is, we cannot recover the **order** of the words.
- We need another indicator vector to reveal the **position** information.

# Position Information III

- To do so, first we choose a maximum length  $T$ .
- No matter how long the document is, we always consider only up to the  $T$ th word.
- Then we can expand the indicator vector in (1) to

$$\mathbf{x}_{i,1:j} \in \{0, 1\}^{|\text{Vocabulary}| \times T}$$

so that

$$\mathbf{x}_{i,1:j} = \begin{matrix} & & & 7 & & & & 25 & & \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ \vdots \end{matrix} & \begin{bmatrix} 0 & \dots & 0 & 0 & 0 & \dots & 0 & 1 & 0 & \dots \\ 0 & \dots & 0 & 1 & 0 & \dots & 0 & 0 & 0 & \dots \\ & & & \vdots & & & & & & \\ & & & \vdots & & & & & & \end{bmatrix} \end{matrix}.$$

# Position Information IV

- Our  $x_{i,1:j}$  is now a matrix. If we insist on that the input is a vector, we can do a conversion by, for example, concatenating all columns.
- With the position information, we can exactly recover our word sequence.

# Output of the Network I

- For output, assume our network can generate

$$\hat{y}_{ij} = f(\theta; x_{i,1:j}) \in [0, 1]^{\text{Vocabulary}}$$

indicating the occurrence probability of each word.

- Consider the example in (2). To predict the next token, we have

$y_{ij}$  is “a”

from the ground truth of the training data.

# Output of the Network II

- Thus we hope that our predicted  $\hat{\mathbf{y}}_{ij}$  satisfies

$$(\hat{\mathbf{y}}_{ij})_{100} \approx 1$$

and

other elements of  $\hat{\mathbf{y}}_{ij} \approx 0$ .

# Autoregressive Settings I

- Up to now what we have is

$$\begin{array}{c} f(\boldsymbol{\theta}; \mathbf{x}_{i,1:j}) \in [0, 1]^{|Vocabulary|} \\ \uparrow \\ \text{network } f(\boldsymbol{\theta}; \mathbf{x}_{i,1:j}) \\ \uparrow \\ \mathbf{x}_{i,1:j} \in \{0, 1\}^{|Vocabulary| \times T} \end{array}$$



# Autoregressive Settings II

- Everything looks good so far. However, as described later, an issue is that we treat

$$\mathbf{x}_{i,1:1}, \mathbf{x}_{i,1:2}, \dots$$

as **independent vectors**.

- We have inputs like

|                      |        |
|----------------------|--------|
| $\mathbf{x}_{i,1:1}$ | I      |
| $\mathbf{x}_{i,1:2}$ | I am   |
| $\mathbf{x}_{i,1:3}$ | I am a |

# Autoregressive Settings III

- They are related. For example,  $\mathbf{x}_{i,1:1}$  is a sub-string of  $\mathbf{x}_{i,1:2}$ , so **we should not treat them as independent vectors**.
- To have efficient training and prediction, we must have a way so that the same operation is not conducted multiple times.
- We know

$\mathbf{x}_{i,1:j}$  includes  $x_{i,1}, x_{i,2}, \dots, x_{i,j}$ .

# Autoregressive Settings IV

- Therefore, we should design  $f(\cdot)$  a way so that the following sequence of calculation can be efficiently done.

$$f(\boldsymbol{\theta}; \mathbf{x}_{i,1:1}), \dots, f(\boldsymbol{\theta}; \mathbf{x}_{i,1:j}), \dots \quad (3)$$

- That is, we hope things are coupled as indicated in Figure 1.

# Autoregressive Settings V

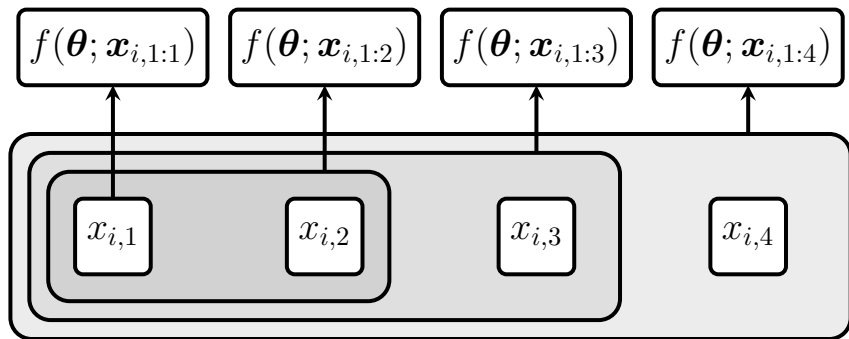


Figure: A sequence of next-token predictions

# Autoregressive Settings VI

- But ensuring that operations used to calculate  $f(\boldsymbol{\theta}; \mathbf{x}_{1:j-1})$  can also be for  $f(\boldsymbol{\theta}; \mathbf{x}_{1:j})$  is not an easy task. Subsequently we introduce an auto-regressive setting for our purpose.

# Summary I

- What we have shown is a high-level illustration of LLM.
- The remaining task is to design a suitable function  $f$ .