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-contextThe next tokenI \rightarrow amI am \rightarrow aI am a \rightarrow machineI am a machine \rightarrow learningI am a machine learning \rightarrow researcher

February 14, 2025

LLM Basics II

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

 $oldsymbol{x}_{i,1:j}$: contexts from 1st to jth words $oldsymbol{y}_{ij}$: (j+1)st word.

• Then we have many

(target, instance)

pairs for training.

• For example, we have

 $oldsymbol{x}_i$: "I am a machine learning researcher" and

LLM Basics III

 $egin{array}{cccc} oldsymbol{x}_{i,1:1} & \mathsf{I} \ oldsymbol{x}_{i,1:2} & \mathsf{I} \ \mathsf{am} \ dots \end{array}$

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

$$\min_{oldsymbol{ heta}} \sum_{i=1}^{D} \sum_{j} \xi(f(oldsymbol{ heta}; oldsymbol{x}_{i,1:j}), oldsymbol{y}_{ij})$$

to obtain the model, where

• heta is the model parameters,

LLM Basics IV

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

$$oldsymbol{y}_{ij}$$
 and $oldsymbol{x}_{i,1:j}$

February 14, 2025

4 / 22

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

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

to ensure that $\boldsymbol{\theta}$ leads to an output as close to \boldsymbol{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,

• automatic differentiation,

then we "believe" that an LLM can be implemented.

- We will explain how we can convert $(\boldsymbol{y}_{ij}, \boldsymbol{x}_{i,1:j})$ to vectors.
- Subsequently, we abuse the notation a bit so that *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:

DocumentI am a machine learning researcherTokenized textIamamachineIearningimage: text imageimage: text imageToken IDs2571002101000

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

February 14, 2025

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

$$\boldsymbol{x}_{i,1:j} \in \{0,1\}^{|\mathsf{Vocabulary}|}$$
 (1)

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

• For example, if

$$oldsymbol{x}_{i,1:2}$$
 is "I am" (2)

February 14, 2025

Position Information II

then we have

$$m{x}_{i,1:2} = [\overbrace{0,\ldots,0,1}^{25 \mathrm{th}}, 0,\ldots,0,1,0,\ldots]^T.$$

- However, from the indicator vector and the Vocabulary, we can extract only the set of words in *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

 $oldsymbol{x}_{i,1:j} \in \{0,1\}^{|\mathsf{Vocabulary}| imes T}$

so that

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

$$\hat{oldsymbol{y}}_{ij} = f(oldsymbol{ heta}; oldsymbol{x}_{i,1:j}) \in [0,1]^{|\mathsf{Vocabulary}|}$$

indicating the occurrence probability of each word.

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

$$oldsymbol{y}_{ij}$$
 is "a"

from the ground truth of the training data.

Output of the Network II

• Thus we hope that our predicted $\hat{y_{ij}}$ satisfies

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

and

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

Autoregressive Settings I

• Up to now what we have is

$$f(oldsymbol{ heta}; oldsymbol{x}_{i,1:j}) \in [0, 1]^{|\mathsf{Vocabulary}|}$$
 \uparrow
network $f(oldsymbol{ heta}; oldsymbol{x}_{i,1:j})$
 \uparrow
 $oldsymbol{x}_{i,1:j} \in \{0, 1\}^{|\mathsf{Vocabulary}| imes T}$

イロト イポト イヨト イヨト 二日

Autoregressive Settings II

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

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

as independent vectors.

• We have inputs like

 $egin{array}{cccc} oldsymbol{x}_{i,1:1} & \mathsf{I} \ oldsymbol{x}_{i,1:2} & \mathsf{I} \ \mathsf{am} \ oldsymbol{x}_{i,1:3} & \mathsf{I} \ \mathsf{am} \ \mathsf{am} \ \mathsf{am} \end{array}$

◆□ ▶ < 圕 ▶ < 壹 ▶ < 壹 ▶</p>
February 14, 2025

Autoregressive Settings III

- They are related. For example, x_{i,1:1} is a sub-string of 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

$$x_{i,1:j}$$
 includes $x_{i,1}, x_{i,2}, \ldots, x_{i,j}$.

February 14, 2025

Autoregressive Settings IV

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

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

February 14, 2025

19/22

• 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(θ; x_{1:j-1}) can also be for f(θ; x_{1:j}) is not an easy task. Subsequently we introduce an auto-regressive setting for our purpose.



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