

Instructions for ACL-2015 Proceedings

Weiming Luo

Student id:mc464633

University of Macau

Department of data science

Macao, China

mc46463@um.edu.mo

Abstract

This modeling project used different languages to perform transfer learning on the English pre-trained model llama2 to allow the model to continue the story in different languages, with good results

1 Introduction

The purpose of this project is to use the pre-training model llama2 for story generation, where the model is given a prompt for the start of the story and the model follows the prompt to continue the story. The model has been pre-trained on the English dataset, and now the model is subjected to transfer learning so that it can continue the story in Chinese and Portuguese.

2 Dataset and training

All the data set are .jsonl format, per line of data set is like

```
{text:"Once upon a time,...."}
```

English, Chinese, Portuguese test data set 500 line, Chinese, Portuguese validation set 500 line, Chinese, Portuguese training set 10000 line. Additional Japanese dataset 400 for training, 10 for validation and testing. And training at free gpu server which support by Kaggle (Kaggle, 2023).

3 Methodology

3.1 Task 1 Decode dataset

The purpose of this section is to code the dataset, and the code for each word is the position of this word in the whole lexicon. Because this project uses pytorch as a framework, it returns the tensor in torch format.

Model.generate() is the function which can continue the story, then generate the encode of the story

```
# Step 1: Encode raw text using tokenizer model. Run tokenization and covert strings into token ids in vocabulary.
tokenized_input = tokenizer.encode(prompt, return_tensors='pt').to(device)
# See the tokenized results.
print(tokenized_input)

tensor([[ 1, 90909, 12081, 263, 931, 20800, 624, 3547, 4562, 750,
         263, 12561, 29880]])
```

Figure 1: prompt encoding

[illegible]

Figure 2: encoding of answer

Then use decoder to decode the story form tensor to English. The principle is to use the index value of each word to find the corresponding word from the original thesaurus, and then convert it into English to print it out.

```
# P 3.11.1 Convert model outputs into raw text.
# Shows how to use the toText() function.
print("I am a ModelOutput object.")
# We only have one data instance, so we directly choose the first row of model output, i.e., output[0][0].
print(bowmodel.toText(output[0][0], skip_special_tokens=True))

# P 3.11.2 Convert model outputs into raw text.
# Shows how to use the toText() function.
print("I am a ModelOutput object.")
# We only have one data instance, so we directly choose the first row of model output, i.e., output[0][0].
print(bowmodel.toText(output[0][0], skip_special_tokens=True))
```

Figure 3: decode of answer

Then change some of the parameters, unlike the last generation, this time we chose to sample randomly and set the temperature to 0.6, indicating that the generated text has some randomness but is still generally skewed towards determinism (the higher the temperature the more random it is, the lower it is the more we tend to favor more probable answers), The top_p and top_k parameters were then adjusted as prompted, which together can determine the diversity and randomness of the sampling process, and when these two parameters were turned up to lower p to lower k, the model generated statements that could withstand more logical scrutiny. However, these few parameter changes didn't cause too many prob-

lems with the model’s logic for generating the articles, but when I tried to change the prompt to a format and emotional coloring that wasn’t quite the same as in the dataset (unlike the once upon a time there was a day with this kind of formatted beginning, and not the happy coloring of a fairy tale), the logic for generating the story started to get messed up. Overall the model does a good job

o.o Stella was hurt because, she had been playing in the garden and she wanted to go inside. Her nanny said, "No, Stella no, you must stay outside." Stella was very sad, she wanted to go inside and have a snack.

Stella's nanny said, "If you stay outside, I will write you a special letter." Stella was so excited. She asked, "What is it?" Nanny said, "It's a surprise!" Stella was so happy she started to run around the garden. She was so excited she forgot to be careful. She ran so fast that she got a fall.

Stella fell down and started to cry. Nanny came to help her and said, "Stella, you must be careful when you run. You need to stay safe when you are going." Stella was so sorry she had done so carelessly. She said, "I'm sorry, Nanny. I will be more careful next time." Nanny smiled and said, "That's good. Now let's go inside and have a snack."

Figure 4: Enter Caption

with English stories, when trying other languages such as Noon, the model doesn't generate along with Chinese, but instead generates English stories that have nothing to do with prompts

有一天，一只狗跑进来了。He was a very happy boy who loved to play outside. One day, he was playing in the park and he saw a big, red ball. He wanted to play with it, so he ran over and grabbed it. But he had never seen a dog before. It was a big, scary dog! The dog barked and growled and started to chase him. He was so scared that he dropped the ball and ran away. The dog chased him all the way home. He was so scared that he didn't even make it. He was so scared that he ran into his bed and hid under the covers. The next day, the boy was still scared. He was too scared to go back to the park. He stayed in his bed all day and all night. He was so scared that he never went to school.

Figure 5: Sample of Chinese prompt

3.2 Task2 Perplexity

Firstly, the model evaluation functions Perplexity of this model is introduced, which in probabilistic terms is the average uncertainty of the language model with respect to the test data. t is the number of tokens in the sequence $P(x_i)$ is the probability that the model will give the t token when it is generated. If perplexity is low, it means that the model predicts the data more accurately.

$$\text{Perplexity} = \exp \left(-\frac{1}{T} \sum_{t=1}^T \log P(x_t \mid x_{<t}) \right)$$

Dataset was test for four different languages and calculated their perplexity values, we found that except for English, which is very low, the perplexity values of Japanese, Chinese and Portuguese are very high, probably because the model is a pre-training model of English, which is a good fit for English grammar, and Portuguese, which is also made up of alphabets, has a higher degree of similarity to English. The perplexity values are also lower than those of Japanese. The perplexity values are lower than those of Japanese and Chinese.

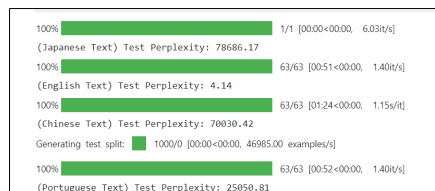


Figure 6: perplexity of four language

3.3 Task 3 Chinese transfer learning

The purpose of this task is to perform transfer learning on the model, and we will use this pre-trained model for English to train Chinese and other languages. Firstly, we used default parameters to train Chinese data set for 8 epochs, At the end, the loss value has converged and looks good, so let's experiment with the model

```
data_file = TEST_FILE
test_dataset = load_dataset('json', data_files={'test': data_file})['test']['text']

results = compute_gpl(model=model, tokenizer=tokenizer, device=device, inputs=test_dataset, batch_size=4)
dataset_results = results['mean_perplexity']
print(f"Test Perplexity: {dataset_results[0].f}")
```

Figure 7: Perplexity of Chinese

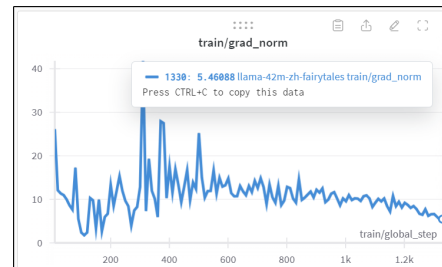
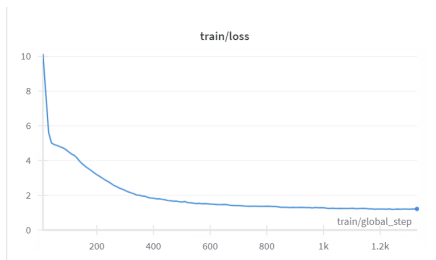
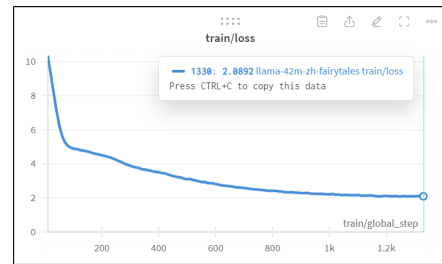
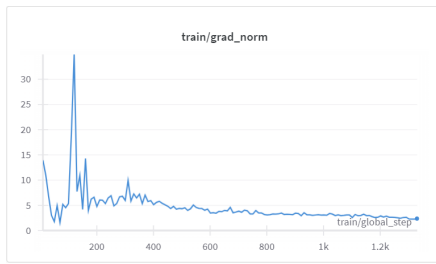
After training the model, Chinese perplexity of the model is test and find that it is reduced to the same level as English, indicating that the model has a full understanding of Chinese. I test the model by giving it a Chinese prompt, and find that the model can already continue the story in Chinese.

```
prompt = "一只从法国海狗嘴里，海狗和章鱼黑黑在一起玩。"  
  
# Decoding hyperparameters  
max_new_tokens = 300  
do_sample = True  
temperature = 0.3  
  
tokenizer_input = tokenizer.encode(prompt, return_tensors="pt").to(device)  
output_ids = model.generate(  
    tokenizer_input,  
    max_new_tokens=max_new_tokens,  
    eos_token_id=1,  
    do_sample=do_sample,  
    temperature=temperature,  
)  
output_text = tokenizer.decode(output_ids[0])  
print(output_text)
```

Figure 8: Enter Caption

We gave a Chinese prompt to test the model, and found that the model can already be used in Chinese to continue the story, we changed the parameters of the generator such as task1, and found that the model’s results gave the expected answers, but in general the logic is not as smooth as in English, probably because the difficulty of the Chinese language is greater than that of the English language.

For the control experiment, the learning rate of the model was tuned to $3e-5$, and the learning rate was lowered to see how it affects the losses



After 8 epochs of training, the model's $grad_{norm}$ is still in a fluctuating state and has not been reduced to 0, and the loss is larger than the first training, which proves that 8 epochs of learning after reducing the learning rate is not enough. In re-testing, the model gives answers that are also illogical.

took a lower learning rate of 3e-5, and 90 epochs of training for such an approach.



3.4 Fonts

And model did not give the answer I want. Therefore there 8 epochs additional training for this model, when epochs attach 16, loss decay to 1.29, as same as learning rate 1e-4. And the Chinese perplexity value of the model has been reduced to 4, and the model has a good understanding of the Chinese language as well

