# 了解 LLaMA-2 模型结构(5)

发表评论 / ChatGPT, GPT, OpenAI

## 9. 转换 tokenizer.model 并保存

前面的章节转换模型的所有权重后，还需要转换 tokenizer.model 为自己需要的格式。

把 meta-llama/Llama-2-7b-chat-hf/ 目录下的 tokenizer.model 拷贝到 newsrc 目录下。

参照 https://github.com/karpathy/llama2.c 项目下的 tokenizer.py 文件，命名为 test09.py，文件保存到 newsrc 目录下：

```python
import os
import struct
import argparse
from typing import List

from sentencepiece import SentencePieceProcessor

class Tokenizer:
    def __init__(self, tokenizer_model=None):
        model_path = tokenizer_model
        assert os.path.isfile(model_path), model_path
        self.sp_model = SentencePieceProcessor(model_file=model_path)
        self.model_path = model_path

        # BOS / EOS token IDs
        self.n_words: int = self.sp_model.vocab_size()
        self.bos_id: int = self.sp_model.bos_id()
        self.eos_id: int = self.sp_model.eos_id()
        self.pad_id: int = self.sp_model.pad_id()
        #print(f"#words: {self.n_words} - BOS ID: {self.bos_id} - EOS ID: {self.eos_id}")
        assert self.sp_model.vocab_size() == self.sp_model.get_piece_size()

    def encode(self, s: str, bos: bool, eos: bool) -> List[int]:
        assert type(s) is str
        t = self.sp_model.encode(s)
        if bos:
            t = [self.bos_id] + t
        if eos:
            t = t + [self.eos_id]
        return t

    def decode(self, t: List[int]) -> str:
        return self.sp_model.decode(t)

    def export(self):

        # get all the tokens (postprocessed) and their scores as floats
        tokens, scores = [], []
        for i in range(self.n_words):

            # decode the token and light postprocessing
            t = self.sp_model.id_to_piece(i)
            s = self.sp_model.get_score(i)
            if i == self.bos_id:
                t = '\n<s>\n'
            elif i == self.eos_id:
```

```python
47              t = '\n</s>\n'
48              t = t.replace('_', ' ') # sentencepiece uses this character as whitespace
49              b = t.encode('utf-8') # bytes of this token, utf-8 encoded
50
51              tokens.append(b)
52              scores.append(s)
53
54          # record the max token length
55          max_token_length = max(len(t) for t in tokens)
56
57          # write to a binary file
58          # the tokenizer.bin file is the same as .model file, but .bin
59          tokenizer_bin = self.model_path.replace('.model', '.bin')
60          with open(tokenizer_bin, 'wb') as f:
61              f.write(struct.pack("I", max_token_length))
62              for bytes, score in zip(tokens, scores):
63                  f.write(struct.pack("fI", score, len(bytes)))
64                  f.write(bytes)
65
66  t = Tokenizer("newsrc/tokenizer.model")
67  t.export()
```

运行 `test09.py`，查看`newsrc` 的文件目录

```
1  ls -l newsrc/tokenizer.*
2  -rwxrwxrwx 1 tony tony 433869 Mar 11 14:24 newsrc/tokenizer.bin
3  -rwxrwxrwx 1 tony tony 499723 Mar 10 23:51 newsrc/tokenizer.model
```

# 10．查看 `tokenizer.bin`

下面给出的例子开始的都是`C/C++`代码，这样更好理解文件里面的内容

参照 `https://github.com/karpathy/llama2.c` 项目下的 `run.c` 文件，命名为 `test01.c`，文件保存到 `newsrc` 目录下：

```c
1  #include <stdio.h>
2  #include <stdlib.h>
3  // The Byte Pair Encoding (BPE) Tokenizer that translates strings <-> tokens
4
5  typedef struct {
6      // 一个字符指针，指向与ID关联的字符串。
7      char *str;
8      // 一个整数，表示与字符串关联的ID。
9      int id;
10 } TokenIndex;
11
12 typedef struct {
13     // 一个指针数组，存储词汇表中每个单词的字符串表示。
14     // 例如，如果 vocab[0] 是 "apple"，那么 vocab[0][0] 就是字符 'a'。
15     char** vocab;
16     // 一个浮点数数组，可能存储与词汇表中的每个单词相关联的分数或权重。
17     float* vocab_scores;
18     // 结构体，存储一个字符串和与其相关联的整数ID
19     TokenIndex *sorted_vocab;
20     // 一个整数，表示词汇表的大小，即vocab和vocab_scores数组中的元素数量。
21     int vocab_size;
22     // 一个无符号整数，表示词汇表中的最大token长度。
23     unsigned int max_token_length;
24     // 一个无符号字符数组，存储所有单字节字符串
25     // 这个数组的大小被设定为512，可能是为了存储ASCII字符表中的所有可能的单字节字符串。
26     unsigned char byte_pieces[512]; // stores all single-byte strings
27 } Tokenizer;
```

```c
28
29  void build_tokenizer(Tokenizer* t, char* tokenizer_path, int vocab_size) {
30      // i should have written the vocab_size into the tokenizer file... sigh
31      // 设置了Tokenizer结构体中vocab_size的值。
32      t->vocab_size = vocab_size;
33      // malloc space to hold the scores and the strings
34      // 为vocab，vocab_scores数组分配了内存，而sorted_vocab被初始化为NULL，表示它将在后续被"懒惰地"初始化
35      t->vocab = (char**)malloc(vocab_size * sizeof(char*));
36      t->vocab_scores = (float*)malloc(vocab_size * sizeof(float));
37      t->sorted_vocab = NULL; // initialized lazily
38      // 初始化byte_pieces数组，该数组存储所有单字节字符串
39      for (int i = 0; i < 256; i++) {
40          t->byte_pieces[i * 2] = (unsigned char)i;
41          t->byte_pieces[i * 2 + 1] = '\0';
42      }
43      // read in the file
44      FILE *file = fopen(tokenizer_path, "rb");
45      if (!file) { fprintf(stderr, "couldn't load %s\n", tokenizer_path); exit(EXIT_FAILURE);
46      // 从文件中读取max_token_length的值
47      if (fread(&t->max_token_length, sizeof(int), 1, file) != 1) { fprintf(stderr, "failed rea
48      int len;
49      // 这个循环块读取每个vocab字符串及其对应的vocab_scores值。
50      for (int i = 0; i < vocab_size; i++) {
51          if (fread(t->vocab_scores + i, sizeof(float), 1, file) != 1) { fprintf(stderr, "faile
52          if (fread(&len, sizeof(int), 1, file) != 1) { fprintf(stderr, "failed read\n"); exit
53          t->vocab[i] = (char *)malloc(len + 1);
54          if (fread(t->vocab[i], len, 1, file) != 1) { fprintf(stderr, "failed read\n"); exit(E
55          t->vocab[i][len] = '\0'; // add the string terminating token
56      }
57      fclose(file);
58  }
59
60  void free_tokenizer(Tokenizer* t) {
61      for (int i = 0; i < t->vocab_size; i++)
62      {
63          free(t->vocab[i]);
64      }
65      if(t->vocab) free(t->vocab);
66      if(t->vocab_scores) free(t->vocab_scores);
67      if(t->sorted_vocab) free(t->sorted_vocab);
68  }
69
70  void print_tokenizer(Tokenizer* t) {
71          printf("vocab = %d\n", t->vocab_size);
72          printf("max_token_length = %d\n", t->max_token_length);
73      for (int i = 0; i < t->vocab_size; i++)
74      {
75          printf("%5d, %12.6lf, (%s)\n", i, t->vocab_scores[i], t->vocab[i]);
76      }
77  }
78
79  int main(int argc, char *argv[]) {
80
81      char *tokenizer_path = "tokenizer.bin";
82
83      // build the Tokenizer via the tokenizer .bin file
84      Tokenizer tokenizer;
85      int vocab_size = 32000; // 从模型文件的 config.json 获取
86      build_tokenizer(&tokenizer, tokenizer_path, vocab_size);
87
88      print_tokenizer(&tokenizer);
89      free_tokenizer(&tokenizer);
90
91      return 0;
92  }
```

编译 `test01.c`

```
1  make test01
2  cc      test01.c   -o test01
```

运行 `test01`

```
1  ./test01 > 1.txt
```

由于输出的内容很多，所以我们把输出重定向到 `1.txt` 文件中，下面是 `1.txt` 文件的开头和结尾部分内容

```
1   vocab = 32000
2   max_token_length = 27
3       0,      0.000000, (<unk>)
4       1,      0.000000, (
5   <s>
6   )
7       2,      0.000000, (
8   </s>
9   )
10      3,      0.000000, (<0x00>)
11      4,      0.000000, (<0x01>)
12      5,      0.000000, (<0x02>)
13      6,      0.000000, (<0x03>)
14      7,      0.000000, (<0x04>)
15      8,      0.000000, (<0x05>)
16      9,      0.000000, (<0x06>)
17     10,      0.000000, (<0x07>)
18     11,      0.000000, (<0x08>)
19     12,      0.000000, (<0x09>)
20  ...
21    259, -1000000000.000000, (   )
22    260,     -1.000000, ( t)
23    261,     -2.000000, (er)
24    262,     -3.000000, (in)
25    263,     -4.000000, ( a)
26    264,     -5.000000, (en)
27    265,     -6.000000, (on)
28    266,     -7.000000, ( th)
29    267,     -8.000000, (es)
30    268, -1000000000.000000, (    )
31    269,    -10.000000, ( s)
32    270,    -11.000000, ( d)
33    271,    -12.000000, (at)
34  ...
35  31985, -31726.000000, (怪)
36  31986, -31727.000000, (联)
37  31987, -31728.000000, (역)
38  31988, -31729.000000, (泰)
39  31989, -31730.000000, (백)
40  31990, -31731.000000, (ȯ)
41  31991, -31732.000000, (げ)
42  31992, -31733.000000, (ベ)
43  31993, -31734.000000, (边)
44  31994, -31735.000000, (还)
45  31995, -31736.000000, (黃)
46  31996, -31737.000000, (왕)
47  31997, -31738.000000, (收)
48  31998, -31739.000000, (弘)
49  31999, -31740.000000, (给)
```

可以看到，`token` 的最大长度为27