

Lab1

2026 年 1 月 21 日

目录

1 基本概念	2
1.1 Tokenizer	2
1.1.1 Word-level tokenizer	2
1.1.2 Byte-level tokenizer	2
1.1.3 Subword tokenizer(BPE)	2
2 问题部分	2
2.1 Problem unicode1	2
2.1.1 Answer	3
2.2 Problem unicode2	3
2.2.1 Answer	4
3 Latex 代码模板	4
4 重点总结	5

1 基本概念

1.1 Tokenizer

1.1.1 Word-level tokenizer

将每个单词定义为一个 token，这样相较于 Byte-level tokenizer 来说，对于一个长句子，会产生更少的 token，但由于需要为每个词都定义一个序号，必然会在某些词并不在定义的字典当中，从而无法解析

1.1.2 Byte-level tokenizer

将每一个字符定义为一个 token，其字典大小为 0-255，这样可以覆盖到每一个可见的字母，但缺点是如果一个单词的长度很长，会使得 token 值特别大，在模型训练的时候会更浪费时间

1.1.3 Subword tokenizer(BPE)

结合 Word-level tokenizer 和 Byte-level tokenizer 的优点，得到了 BPE 这种压缩算法，即 Byte-pair encoding

对于一个句子，我们的字典最初就是 Byte-level tokenizer 的字典，序号是 0-255，然后遍历整个句子，每次将相邻的两个单词组成一个 pair，然后记录该 pair 出现的次数，遍历完以后，将出现次数最高的 pair 定义为一个新的字典序号，并加入到字典当中，然后将该 pair 看成是一个 word，持续遍历整个句子，直到合并完成。直接做的缺点有每次都需要完整遍历语料库，计算量非常大，其次，如果是相邻的标点符号不一样，也会生成一个新的 ID，尽管标点前面的单词是一样的

在最初提出的 BPE 中，字典里存在一些 **Special tokens**，比如 `<|endoftext|>`，其用途是告诉程序什么时候停止生成 token，**special tokens** 不会被分割成多个 token，在该实验的 BPE 中不考虑 special token

给定 `low`、`lower`、`widest`、`newest`，根据 BPE 的步骤进行分词

<code>low</code> : 5	<code>[l,ow]</code>	<code>lower</code> : 2	<code>[[l,ow],r]</code>	<code>widest</code> : 3	<code>[[wid,est]]</code>	<code>newest</code> : 6	<code>[[ne,newt]]</code>
<code>(l,o): 7</code>	<code>(l,o): 7</code>	<code>(l,o): 7</code>	<code>(l,ow): 7</code>	<code>(low,e): 2</code>	<code>(low,e): 2</code>	<code>(low,e): 2</code>	<code>(low,e): 2</code>
<code>(o,w): 7</code>	<code>(o,w): 7</code>	<code>(o,w): 7</code>	<code>(ow,e): 2</code>	<code>(e,r): 2</code>	<code>(e,r): 2</code>	<code>(e,r): 2</code>	<code>(low,e): 2</code>
<code>(w,e): 6</code>	<code>(w,e): 8</code>	<code>(w,e): 2</code>	<code>(e,v): 2</code>	<code>(w,i): 3</code>	<code>(w,i): 3</code>	<code>(w,i): 3</code>	<code>(e,r): 2</code>
<code>(e,r): 2</code>	<code>(e,r): 2</code>	<code>(e,r): 2</code>	<code>(w,i): 3</code>	<code>(w,i): 3</code>	<code>(i,d): 3</code>	<code>(i,d): 3</code>	<code>(e,r): 2</code>
<code>(w,i): 3</code>	<code>(w,i): 3</code>	<code>(w,i): 3</code>	<code>(i,d): 3</code>	<code>(i,d): 3</code>	<code>(d,est): 3</code>	<code>(d,est): 3</code>	<code>(e,r): 2</code>
<code>(i,d): 3</code>	<code>(i,d): 3</code>	<code>(i,d): 3</code>	<code>(d,est): 3</code>	<code>(n,e): 6</code>	<code>(n,e): 6</code>	<code>(n,e): 6</code>	<code>(low,e): 2</code>
<code>(d,e): 3</code>	<code>(d,e): 3</code>	<code>(d,e): 3</code>	<code>(d,est): 3</code>	<code>(e,w): 6</code>	<code>(e,w): 6</code>	<code>(e,w): 6</code>	<code>(low,r): 2</code>
<code>(e,s): 9</code>	<code>(e,s): 9</code>	<code>(e,s): 9</code>	<code>(n,e): 6</code>	<code>(e,w): 6</code>	<code>(e,w): 6</code>	<code>(e,w): 6</code>	
<code>(n,s): 9</code>	<code>(n,s): 9</code>	<code>(n,s): 9</code>	<code>(n,e): 6</code>	<code>(n,e): 6</code>	<code>(n,e): 6</code>	<code>(n,e): 6</code>	
<code>(n,e): 6</code>	<code>(n,e): 6</code>	<code>(n,e): 6</code>	<code>(e,w): 6</code>	<code>(e,w): 6</code>	<code>(e,w): 6</code>	<code>(e,w): 6</code>	
<code>(e,w): 6</code>			<code>(w,est): 6</code>	<code>(w,est): 6</code>	<code>(w,est): 6</code>	<code>(w,est): 6</code>	

2 问题部分

2.1 Problem unicode1

Problem (unicode1): Understanding Unicode (1 point)

- (a) What Unicode character does `chr(0)` return?

Deliverable: A one-sentence response.

- (b) How does this character's string representation (`__repr__()`) differ from its printed representation?

Deliverable: A one-sentence response.

- (c) What happens when this character occurs in text? It may be helpful to play around with the following in your Python interpreter and see if it matches your expectations:

```
1 >>> chr(0)
2 >>> print(chr(0))
3 >>> "this is a test" + chr(0) + "string"
4 >>> print("this is a test" + chr(0) + "string")
```

Deliverable: A one-sentence response.

2.1.1 Answer

- (a) 在终端中使用 python 环境进行实现

```
|>>> chr(0)
|'\x00'
>>> □
```

- (b) `chr(0)` 是空字符，直接在命令行输入命令或调用 `__repr__()` 时会用转义字符来表示，对于非打印字符来说直接打印是不可见的

```
|>>> print(chr(0))
>>> □
```

- (c) `print()` 是位于用户视图，直接在交互式命令行输入命令或调用 `__repr__()` 是位于开发者视图，因此一个可见空字符，一个不可见空字符

```
|>>> "this is a test" + chr(0) + "string"
|'this is a test\x00string'
|>>> print("this is a test" + chr(0) + "string")
|this is a teststring
```

2.2 Problem unicode2

Problem (unicode2): Unicode Encodings (3 points)

- (a) What are some reasons to prefer training our tokenizer on UTF-8 encoded bytes, rather than UTF-16 or UTF-32? It may be helpful to compare the output of these encodings for various

input strings.

Deliverable: A one-to-two sentence response.

- (b) Consider the following (incorrect) function, which is intended to decode a UTF-8 byte string into a Unicode string. Why is this function incorrect? Provide an example of an input byte string that yields incorrect results.

```

1 def decode_utf8_bytes_to_str_wrong(bytestring: bytes):
2     return "".join([bytes([b]).decode("utf-8") for b in bytestring])
3
4 >>> decode_utf8_bytes_to_str_wrong("hello".encode("utf-8"))
5 'hhello'
```

Deliverable: An example input byte string for which `decode_utf8_bytes_to_str_wrong` produces incorrect output, with a one-sentence explanation of why the function is incorrect.

- (c) Give a two byte sequence that does not decode to any Unicode character(s).

Deliverable: An example, with a one-sentence explanation.

2.2.1 Answer

- (a) UTF-8 相较于 UTF-16 和 UTF-32 有以下几个优点：

- UTF-8 完全兼容 ASCII 码，而 UTF-16 和 UTF-32 处理 ASCII 字符时会插入大量空字节
- UTF-8 是变长编码，通常是 1 到 4 个字节不等，对比 UTF-16 和 UTF-32 更节省空间
- UTF-8 是以单字节为单位处理的，不存在大端序和小端序转换的问题，而 UTF-16 和 UTF-32 必须处理字节序问题，否则会出现乱码
- 如果文件在传输过程中丢失了一个字节或受损，UTF-8 能根据后续的字节特征快速找到下一个字符的起始位置

- (b) 这段代码会出错，因为在函数当中处理 `bytestring` 的方式是一个字节一个字节的解码，但 UTF-8 是变长编码，如果输入的是多个字节组成的字符会导致 **UnicodeDecodeError**

```

1 def decode_utf8_bytes_to_str_wrong(bytestring: bytes):
2     return "".join([bytes([b]).decode("utf-8") for b in bytestring])
3
4 if __name__ == "__main__":
5     print(decode_utf8_bytes_to_str_wrong("hello,你好".encode("utf-8")))

UnicodeDecodeError: 'utf-8' codec can't decode byte 0xe4 in position 0: unexpected end of data
```

- (c) 0XC0 0XAF，在现代解码器中任何以 0XC0、0XC1 开头的字节序列都会被抛出 **UnicodeDecodeError**

行内公式示例： $E = mc^2$ 。

3 Latex 代码模板

在笔记中插入代码是非常常见的需求，特别是 CS 专业的笔记。

定义 3.1：梯度下降

梯度下降（Gradient Descent）是一种用于优化函数的迭代算法。为了找到函数的局部最小值，我们需要向当前点梯度的**反方向**迈步。

定理 3.1：泰勒公式

如果函数 $f(x)$ 在点 x_0 处具有 n 阶导数，则有：

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(x_0)}{n!} (x - x_0)^n \quad (1)$$

```
1 def hello_world():
2     print("Hello, LaTeX!")
3     return True
```

Listing 1: Python 示例代码

4 重点总结

注意

考试重点：请务必记住上述定理的适用条件，不要混淆！