## 从零实现大语言模型

Lec1: BPE 分词系统的算法设计与实现

Penghao Kuang

SIST, ShanghaiTech ACM, DataTech, Geekpie

2025.10.19

算法优化 2: 文本编码

#### Overview

算法优化 1: 词表训练中的合并算法

设计动机

设计动机

核心算法

算法优化 1: 词表训练中的合并算法

算法优化 2: 文本编码

实操演示

000000

# PART 1 设计动机

## 大语言模型能否直接识别自然语言文本?

Hello helo. I'm ······

- 计算机只能识别处理数字编码。无法直接处理自然语言文本!
- 如何设计一种处理方法,使得自然语言文本能无损地转化成数字编码?

## 核心思路: 按词表映射

#### 假设拥有一个完整的词表:

- 词表中每一条目包含: 词汇、词汇对应的数码编号
- 对照词表,将文本逐词映射为数码序列

#### 例如词表按整词分割规则制定:

```
"Hello":0
" ":1
"helo":2
",":3
"I'm":4
```

则原文本 Hello helo, I'm 的映射结果为: [0,1,2,3,1,4]

## 核心思路: 按词表映射

#### 问题:

- 自然语言非常复杂且语种众多。如何制作完整的词表?
- 词表总大小是极为恐怖的
- 如何正确分词编码?
  - "United Kingdoms" 作为整词还是两个独立词?
  - 中文等无明显分词界限的语言如何分词?("我一把把把把住了")

整词分割规则显然有明显的局限性

## 尝试思路 2: utf-8 编码分割

#### 我们存在一种现成的文本编码方式:

```
>>> print("h".encode("utf-8"))
h'h'
>>> print("你".encode("utf-8"))
b'\xe4\xbd\xa0'
>>> print("ô".encode("utf-8"))
b'\xc3\xb4'
>>> test string="hello,你好"
>>> utf_encoded=test_string.encode("utf-8")
>>> print(utf_encoded)
b'hello.\xe4\xbd\xa0\xe5\xa5\xbd'
>>> list(utf encoded)
[104, 101, 108, 108, 111, 44, 228, 189, 160, 229, 165, 189]
```

## 尝试思路 2: utf-8 编码分割

#### 我们只需 256 种数字编码, 就可表示所有类型的字符!

```
>>> test_string="hello,你好"
>>> utf_encoded=test_string.encode("utf-8")
>>> print(utf_encoded)
b'hello,\xe4\xbd\xa0\xe5\xa5\xbd'
>>> list(utf_encoded)
[104, 101, 108, 108, 111, 44, 228, 189, 160, 229, 165, 189]
```

我们可以直接按照现有的 utf-8 编码规范, 将文本转化成 utf-8 编码序列!

## 尝试思路 2: utf-8 编码分割

#### 是否存在问题?

设计动机

- 相比整词分割, 文本序列长度将显著膨胀
- 字符的信息密度显著下降,不利于模型捕捉文本信息

能否找到介于"整词分割"与"utf-8编码"之间的折中方案?

⇒BPE 分词算法将有效解决词表制作与分词规则的上述问题!

# PART 2 核心算法

- 1. 初始条件与 utf-8 编码完全一致: 拥有一个大小为 256 的映射词表
- 2. 每次选取语料库中相邻频率最高的字符对进行合并,形成一个全新词条

#### 对于原文本 Hello helo, I'm:

- 1. 建立大小为 256 的原词表 (编号 0-255)
- 2. 映射为序列 [72, 101, 108, 108, 111, 32, 104, 101, 108, 111, 44, 32, 73, 39, 109]
- 3. 计算字符对相邻频率: (e,l):2 (l,o):2 (H,e):1 (l,l):1……
- 4. 取最大频率中字典序最大的词条 (I,o): 所有相邻的 (I,o) 合并成新字符"lo"
- 5. 词表增加词条: 1o:256
- 6. 语料库变为:H,e,1,lo, ,h,e,lo, ,I,',m
- 7. 下一轮合并的相邻频率变为: (e,l):1 (e,lo):1 (lo, ):2……

#### 但在实际工程中,我们还需要考虑更多的细节

- 语料库可能由大量无关的文本拼接而成,由"<|endoftext|>" 符号分开
- 我们倾向于将整个文本粗糙地分成一系列"整词", 再运行 BPE 算法
  - 例如原文本期望分割成: "Hello" " helo" "," " I" "'m"
  - 这样可以有效降低空格、标点、特殊符号对词结构本身的影响
  - 并且避免无意义的跨词合并: 若出现了一个合成结果"Hello helo,"——这是否有实际意义?

u don't have to be scared of the loud dog, I'll protect you kept him safe. The mole had found his best friend.
<|endoftext|>
Once upon a time, in a warm and sunny place, there was a bi
Tom asked his friend, Sam, to help him search for the ball.
Sam and Tom went close to the pit. They were scared, but th
They went into the pit to search. It was dark and scary. Th
but no one could hear them. They were sad and scared, and t
<|endoftext|>

- 原始词表大小调整为 257:256 个基本映射词条 + ("<|endoftext|>": 256) 词条
- 利用正则表达式进行"整词"粗略分割,以此限定合并的极限这就是我们的"Pre-tokenization"(预词元分割)步骤。

## 词表制作: Pre-tokenization

#### 如何对一个语料库进行预词元分割?

```
intj intj intj intj tech tech<|endoftext|>
```

Héllò hôw <|endoftext|><|endoftext|> are ü?

- 1. 先以 < |endoftext|> 为界,拆成若干个文本单位
- 2. 在每个文本单位中,用正则表达式进行"整词"拆分

### 词表制作: Pre-tokenization

```
intj intj intj intj tech tech<|endoftext|>
Héllò hôw <|endoftext|><|endoftext|> are ü?
```

#### 分割结果:

```
chunk a text:b'intj intj intj intj intj tech tech'
chunk a text:b'<|endoftext|>'
chunk a text:b'\n\nH\xc3\xa911\xc3\xb2 h\xc3\xb4w '
chunk a text:b'<|endoftext|>'
chunk a text:b'<|endoftext|>'
chunk a text:b' are \xc3\xbc?'
```

### 词表制作: Pre-tokenization

#### 采用正则表达式分割。GPT-2 的分词标准为:

#### 部分文本的切割效果:

设计动机

```
chunk a text:b'intj intj intj intj intj tech tech'
>>>>split a token:b'intj'
>>>>split a token:b' intj'
>>>>split a token:b' intj'
>>>>split a token:b' intj'
>>>>split a token:b' intj'
>>>>split a token:b' tech'
>>>>split a token:b' tech'
chunk a text:b'<|endoftext|>'
>>>>split a token:b'<|endoftext|>'
```

代码讲解: Pre-tokenization

## 词表制作: 合并

#### 一种朴素的合并方法是:

- 将每个"整词"全部打散为单个 utf-8 编码组成的序列
- 跑遍整个语料库的所有打散序列,统计相邻词对频率
- 按照 BPE 算法,选择最高频词对合并
- 再次跑遍整个语料库的所有打散序列,更改需要合并的打散字符
- 一次合并完成,循环进行后续合并

## 词表制作: 合并

#### 最后得出的词表示例:

```
423: b'ive', 424: b' al', 425: b"'s", 426: b' su', 427: b'el', ......
```

418: b' by', 419: b'os', 420: b'ant', 421: b' G', 422: b' from',

5569: b' turning', 5570: b' Please', 5571: b' Chief', 5572: b' quant'

•••

9801: b' Elizabeth', 9802: b'igate', 9803: b'TR', 9804: b' cognitive'

简单词根与小词 → 高频整词 → 专有名词、次高频词根、次高频词等 往往只需要万级的词表大小,就可以高效表示出大多数的整词及几乎所有的词根 ⇒ 高频词可直接整体编码、低频词可按词根编码

## 二、文本编码

- 在形成词表的过程中同步记录合并的词对。如 (b' ', b't'), (b' ', b'a'), (b'h', b'e'), (b'i', b'n'), (b'r', b'e'), (b' t', b'he'), (b'o', b'n'), (b'e', b'r')
- 对于待编码的文本: 按照合并列表的顺序进行合并:
- 将语料库中所有相邻的 (b' ', b't') 合并为 b' t'
- 将语料库中所有相邻的 (b' ', b'a') 合并为 b' a'
- 将语料库中所有相邻的 (b'h', b'e') 合并为 b'he'
- 将语料库中所有相邻的 (b'i', b'n') 合并为 b'in'

由此可将最初的一长串 utf-8 编号不断合并,合并为一系列可表征整词、有意义词根的词条编号!

## 耗时问题

- 对于极小规模的语料库,不需要考虑耗时问题
- 正规的词表制作中, 语料库往往是几百 MB 规模
- 每一次合并都需跑遍整个语料库进行统计和更改……
- 对于文本编码的合并操作同样需要大量遍历整个语料库
- 上万大小的词表将对应上万次的语料库遍历

以 130KB 大小的语料库为例,上述算法的实际运行时长约为 7 秒。对于几百 MB 级别的语料库,这种算法的时间开销是不可承受的!

## PART 3

算法优化 1: 词表训练中的合并算法

•000000000

## 算法优化 1: 词表训练中的合并算法

#### 原算法的耗时问题在干:

- 每次合并、都对整个语料库中(包括需要合并和无需合并)的内容进行遍历。
- 每次合并都需要生成全新的词对计数表(许多词对相比合并前并没有变化) 所以我们期望设计出一种算法:
  - 能够维护语料库中的词对信息,使得每次合并都能精准更新产生变化的词对
  - 能够维护各词对的计数,使得每次合并只需精准更新计数表的少量条目

算法优化 2: 文本编码

## 需要准备的数据结构

#### 词对信息的维护:

- 为所有字符设置"位置编号"(token\_idx)
- 建立字典, key 为词对, value 为该词对的所有位置: 词对位置设定为第一字符的位置编号(pair\_positions)
- 设置 prev、next 字典, key 为位置编号, value 为前、后的位置编号。可能为 None。
- 设置词对计数器字典。key 为词对, value 为该词对的总出现次数。 (pair\_counter)
- 设置优先队列维护词对计数信息。(heap)

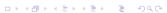
## token idx 的预处理

```
for idx in list(token dict.keys()):
print(f"Token idx: {idx}, Token: {token dict[idx]},
    Char: {int2utf[token dict[idx]]}")
```

```
Token idx: 0, Token: 105, Char: b'i
                                        Token idx: 30. Token: 116. Char: b't'
Token idx: 1, Token: 110, Char: b'n
                                        Token idx: 31, Token: 101, Char: b'e'
Token idx: 2. Token: 116. Char: b't
                                        Token idx: 32, Token: 99, Char: b'c'
Token idx: 3. Token: 106. Char: b'i
                                        Token idx: 33. Token: 104. Char: b'h'
Token idx: 4, Token: 32, Char: b'
Token idx: 5, Token: 105, Char: b'i'
                                        Token idx: 47. Token: 10. Char: b'\n'
                                        Token idx: 48, Token: 10, Char: b'\n'
Token idx: 6, Token: 110, Char: b'n
                                        Token idx: 49. Token: 72. Char: b'H'
Token idx: 7. Token: 116. Char: b't'
Token idx: 8. Token: 106. Char: b'i
                                        Token idx: 50, Token: 195, Char: b'\xc3
                                        Token idx: 51, Token: 169, Char: b'\xa9
Token idx: 9. Token: 32. Char: b' '
Token idx: 10, Token: 105, Char: b'i'
                                        Token idx: 52, Token: 108, Char: b'l'
Token idx: 11, Token: 110, Char: b'n'
                                        Token idx: 53. Token: 108. Char: b'l'
Token idx: 12. Token: 116. Char: b't'
                                        Token idx: 54. Token: 195. Char: b'\xc3
Token idx: 13. Token: 186. Char: b'i'
                                        Token idx: 55, Token: 178, Char: b'\xb2
Token idx: 14, Token: 32, Char: b' '
                                        Token idx: 56. Token: 32. Char: b' '
Token idx: 15. Token: 105. Char: b'i
                                        Token idx: 57, Token: 104, Char: b'h'
Token idx: 16, Token: 110, Char: b'n'
                                        Token idx: 58, Token: 195, Char: b'\xc3
Token idx: 17. Token: 116. Char: b't
                                        Token idx: 59, Token: 180, Char: b'\xb4
Token idx: 18. Token: 186. Char: b'i
                                        Token idx: 60, Token: 119, Char: b'w'
Token idx: 19, Token: 32, Char: b' '
                                        Token idx: 61, Token: 32, Char: b' '
Token idx: 20. Token: 105. Char: b'i
                                        Token idx: 88. Token: 32. Char: b' '
Token idx: 21, Token: 110, Char: b'n
                                        Token idx: 89. Token: 97. Char: b'a'
Token idx: 22. Token: 116. Char: b't
                                        Token idx: 90. Token: 114. Char: b'r'
Token idx: 23. Token: 186. Char: b'i
                                        Token idx: 91, Token: 101, Char: b'e'
Token idx: 24, Token: 32, Char: b' '
                                        Token idx: 92. Token: 32. Char: b' '
Token idx: 25. Token: 116. Char: b't'
                                        Token idx: 93, Token: 195, Char: b'\xc3
Token idx: 26, Token: 101, Char: b'e
                                        Token idy: 94. Token: 188. Char: h'\vhc'
Token idx: 27. Token: 99. Char: b'c'
                                       Token idx: 95. Token: 63. Char: b'?'
Token idx: 28. Token: 184. Char: b'h'
                                        Token idx: 96. Token: 32. Char: b' '
Token idx: 29, Token: 32, Char: b'
```

## pair\_positions 的预处理

```
Pair: (105, 110), Char Pair:((b'i', b'n')), Positions: [0, 5, 10, 15, 20]
Pair: (110, 116), Char Pair: ((b'n', b't')), Positions: [1, 6, 11, 16, 21]
Pair: (116, 106), Char_Pair:((b't', b'j')), Positions: [2, 7, 12, 17, 22]
Pair: (32, 105), Char Pair:((b' ', b'i')), Positions: [4, 9, 14, 19]
Pair: (32, 116), Char Pair: ((b' ', b't')), Positions: [24, 29]
Pair: (116, 101), Char Pair:((b't', b'e')), Positions: [25, 30]
Pair: (101, 99), Char_Pair:((b'e', b'c')), Positions: [26, 31]
Pair: (99, 104), Char Pair:((b'c', b'h')), Positions: [27, 32]
Pair: (72, 195), Char_Pair:((b'H', b'\xc3')), Positions: [49]
Pair: (195, 169), Char_Pair:((b'\xc3', b'\xa9')), Positions: [50]
Pair: (169, 108), Char Pair:((b'\xa9', b'l')), Positions: [51]
Pair: (108, 108), Char Pair: ((b'l', b'l')), Positions: [52]
Pair: (108, 195), Char_Pair:((b'l', b'\xc3')), Positions: [53]
Pair: (195, 178), Char Pair:((b'\xc3', b'\xb2')), Positions: [54]
Pair: (32, 104), Char Pair: ((b' ', b'h')), Positions: [56]
Pair: (104, 195), Char Pair: ((b'h', b'\xc3')), Positions: [57]
Pair: (195, 180), Char_Pair:((b'\xc3', b'\xb4')), Positions: [58]
Pair: (180, 119), Char Pair:((b'\xb4', b'w')), Positions: [59]
Pair: (32, 97), Char Pair:((b' ', b'a')), Positions: [88]
Pair: (97, 114), Char Pair:((b'a', b'r')), Positions: [89]
Pair: (114, 101), Char Pair:((b'r', b'e')), Positions: [90]
Pair: (32, 195), Char Pair:((b' ', b'\xc3')), Positions: [92]
Pair: (195, 188), Char Pair:((b'\xc3', b'\xbc')), Positions: [93]
```



## pair counter 的预处理

```
for pair in list(global pair counter.keys()):
print(f"Pair: {pair}, Char Pair:({int2utf[pair[0]],int2utf[pair[1]]}),
    Counted: {global pair counter[pair]}")
```

```
Pair: (105, 110), Char Pair: ((b'i', b'n')), Counted: 5
Pair: (110, 116), Char Pair: ((b'n', b't')), Counted: 5
Pair: (116, 106), Char Pair: ((b't', b'j')), Counted: 5
Pair: (32, 105), Char_Pair:((b' ', b'i')), Counted: 4
Pair: (32, 116), Char_Pair:((b' ', b't')), Counted: 2
Pair: (116, 101), Char Pair: ((b't', b'e')), Counted: 2
Pair: (101, 99), Char Pair:((b'e', b'c')), Counted: 2
Pair: (99, 104), Char Pair: ((b'c', b'h')), Counted: 2
Pair: (72, 195), Char Pair:((b'H', b'\xc3')), Counted: 1
Pair: (195, 169), Char Pair: ((b'\xc3', b'\xa9')), Counted: 1
Pair: (169, 108), Char Pair: ((b'\xa9', b'l')), Counted: 1
Pair: (108, 108), Char Pair: ((b'l', b'l')), Counted: 1
Pair: (108, 195), Char Pair:((b'l', b'\xc3')), Counted: 1
Pair: (195, 178), Char Pair:((b'\xc3', b'\xb2')), Counted: 1
Pair: (32, 104), Char Pair: ((b' ', b'h')), Counted: 1
Pair: (104, 195), Char Pair: ((b'h', b'\xc3')), Counted: 1
Pair: (195, 180), Char Pair: ((b'\xc3', b'\xb4')), Counted: 1
Pair: (180, 119), Char Pair:((b'\xb4', b'w')), Counted: 1
Pair: (32, 97), Char_Pair:((b' ', b'a')), Counted: 1
Pair: (97, 114), Char Pair: ((b'a', b'r')), Counted: 1
Pair: (114, 101), Char_Pair:((b'r', b'e')), Counted: 1
Pair: (32, 195), Char Pair:((b' ', b'\xc3')), Counted: 1
```

## next 的预处理 (prev 同理)

```
for idx in list(global next.keys()):
print(f"Idx: {idx}, Char: {int2utf[global token dict[idx]]},
    Next idx: {global next[idx]}")
```

```
Idx: 0, Char: b'i', Next idx: 1
                                      Idx: 30, Char: b't', Next idx: 31
Idx: 1. Char: b'n'. Next idx: 2
                                      Idx: 31, Char: b'e', Next idx: 32
Idx: 2, Char: b't', Next idx: 3
                                      Tdx: 32. Char: b'c'. Next idx: 33
Idx: 3. Char: b'i'. Next idx: None
                                      Idx: 33. Char: b'h'. Next idx: None
Idx: 4, Char: b' ', Next idx: 5
                                      Idx: 47. Char: b'\n'. Next idx: None
Idx: 5, Char: b'i', Next idx: 6
                                      Idx: 48, Char: b'\n', Next idx: None
Idx: 6. Char: b'n'. Next idx: 7
                                      Idx: 49. Char: b'H'. Next idx: 50
Idx: 7, Char: b't', Next idx: 8
                                      Idx: 50, Char: b'\xc3', Next idx: 51
Idx: 8, Char: b'j', Next idx: None
                                      Idx: 51. Char: b'\xa9'. Next idx: 52
Idx: 9, Char: b' ', Next idx: 10
                                      Idx: 52, Char: b'l', Next idx: 53
Idx: 10, Char: b'i', Next idx: 11
                                      Idx: 53. Char: b'1'. Next idx: 54
Idx: 11. Char: b'n'. Next idx: 12
                                      Idx: 54, Char: b'\xc3', Next idx: 55
Idx: 12, Char: b't', Next idx: 13
                                      Idx: 55. Char: b'\xb2'. Next idx: None
Idx: 13. Char: b'i'. Next idx: None
Idx: 14, Char: b' ', Next idx: 15
                                      Idx: 56, Char: b' '. Next idx: 57
Idx: 15. Char: b'i'. Next idx: 16
                                      Idx: 57, Char: b'h', Next idx: 58
Idx: 16. Char: b'n'. Next idx: 17
                                      Idx: 58, Char: b'\xc3', Next idx: 59
Idx: 17. Char: b't'. Next idx: 18
                                      Idx: 59, Char: b'\xb4', Next idx: 60
Idx: 18. Char: b'i'. Next idx: None
                                      Idx: 60. Char: b'w'. Next idx: None
Idx: 19, Char: b' ', Next idx: 20
                                      Idx: 61, Char: b' ', Next idx: None
Idx: 20. Char: b'i'. Next idx: 21
                                      Idx: 88, Char: b' '. Next idx: 89
Idx: 21, Char: b'n', Next idx: 22
                                      Idx: 89. Char: b'a'. Next idx: 90
Idx: 22, Char: b't', Next idx: 23
                                      Idx: 90. Char: b'r'. Next idx: 91
Idx: 23. Char: b'i'. Next idx: None
                                      Idx: 91. Char: b'e'. Next idx: None
Idx: 24, Char: b' ', Next idx: 25
                                      Idx: 92, Char: b' ', Next idx: 93
Idx: 25. Char: b't'. Next idx: 26
                                      Idx: 93. Char: b'\xc3'. Next idx: 94
Idx: 26, Char: b'e', Next idx: 27
Idx: 27. Char: b'c'. Next idx: 28
                                      Idx: 94. Char: b'\xbc'. Next idx: None
Idx: 28. Char: b'h'. Next idx: None
                                      Idx: 95. Char: b'?'. Next idx: None
                                      Idx: 96. Char: b' '. Next idx: None
```



## 每次合并的更改过程

- 选定需要合并的词对 (pair0,pair1)
- 遍历该词对对应的所有位置
  - 确认字符无误
  - 获取 pair0 的前一字符 prev、pair1 的后一字符 next
  - prev、pair0、pair1、next、new 相关词对的计数器更新
  - pair0、pair1 对应位置编号的字符更新
  - 相关词对位置编号列表更新(只增不删)
  - prev、next 更新
- 优先队列更新(采用 lazy update)

## 实际算法运行

- 选择计数最大的 (b't',b'j'), 合并为 new=b'tj'
- 对于 pair\_positions 中 (b't',b'j') 对应的所有位置 (所有 pairs):
  - 该位置 b't' 前一个字符 prev\_token、b'j' 后一个字符 next\_token
  - (prev\_token,b't')(prev\_token,new)(b'j',next\_token)(new,next\_token) 计数器更新
  - b'j' 不复存在, token\_idx 该下标改为 None; b't' 处改为 b'tj'
  - (prev\_token,new)(new,next\_token) 的 pair\_positions 列表增加元素(lazy update)
  - b'tj' 下标的 next 改为 next\_token 下标; next\_token 下标改为 b'tj' 下标
- 优先队列 heap: 将所有受影响的 pairs 及计数插入(先前数据不删除)

代码讲解: 合并算法优化

## PART 4 算法优化 2: 文本编码

## 逆向思维

- 预处理中经过了 Pre-tokenization: 不同整词之间互不干扰
- 原本编码算法改成:逐个对单一"整词"进行编码,正确性不受影响
- 对于每个"整词": 枚举所有的相邻词对
- 如果该词对可合并(在合并列表中),查询其在合并列表中的位置
- 位置最靠前的一定被最先合并

## 实际算法运行

```
original token: tech
pair:(b' ', b't'),number:0
pair:(b't', b'e'),number:465
pair:(b'e', b'c'),number:505
pair:(b'c', b'h'),number:98
finish a merge, token list:[b' t', b'e', b'c', b'h']
pair:(b' t', b'e'),number:287
pair:(b'e', b'c'),number:505
pair:(b'c', b'h'),number:98
finish a merge, token list: [b' t', b'e', b'ch']
pair:(b' t', b'e'), number:287
pair:(b'e', b'ch'),number:2614
finish a merge, token list: [b' te', b'ch']
pair:(b' te', b'ch').number:5739
finish a merge, token list: [b' tech']
```

代码讲解: 文本编码优化

# PART 5 实操演示

#### 我们将实操演示:

设计动机

- 21MB 语料库上 1000 大小词表的训练
- 32000 大小词表上的文本编码与解码

# 谢谢!