

熵计算实验报告

陈江昊 人工智能学院 202228014628016 [代码地址](#)

实验介绍

1. 利用爬虫工具从互联网上收集大量的中英文文本数据，并对收集到的数据进行清洗；
2. 设计算法并编程实现在收集文本数据上中文汉字和英文字母的概率和熵的计算；
3. 改变文本数据规模，重新计算中文汉字和英文字母的概率和熵，并分析计算结果。

实验数据收集

中文数据

根据URL以及HTML标签的规律，对百度百科不同词条下网页中的文本数据进行爬取。百度百科的URL形式为 "https://baike.baidu.com/view/" + str(i) + ".htm"，规律性强，从而能够方便的获取大量的URL链接。此外，利用 BeautifulSoup 库对网页HTML源码进行分析，从中提取中文文本数据。具体代码实现如下：

```
1 def get_text_ch(url, num, save_path, save_freq=100):
2     os.makedirs(save_path, exist_ok=True)
3     dic_ch = {}
4     for idx in tqdm(range(1, num + 1, 1)):
5         try:
6             r = requests.get(url + str(idx) + '.htm',
7                               headers=headers, proxies=proxies)
8             # r.encoding = r.apparent_encoding
9             soup = BeautifulSoup(r.text, 'html.parser')
10            text = ''
11            for x in soup.find_all('div', {'class': 'para'}):
12                text += x.text
13            dic_ch[url + str(idx) + '.htm'] = text
14            r.close()
15        except:
16            pass
17
18        if idx % save_freq == 0:
19            u = list(dic_ch.keys())
20            t = list(dic_ch.values())
21            result = pd.DataFrame()
22            result["url"] = u
23            result["text"] = t
24            result.to_excel(os.path.join(save_path,
25                                         "./data{}.xlsx".format(idx)))
26            dic_ch = {}
```

英文数据

由于未能找到URL较为规律的英文网站，我们决定对 Wikipedia 上的数据进行爬取，并以 Wikipedia 上某个词条的网页为种子，寻找到该网页中所有链接到 Wikipedia 其他词条的超链接，然后对找到的新链接再进行数据爬取，以此类推。具体代码实现如下：

```
1 def get_text_en(url, num, keyword, save_path, save_feq=100):
2     os.makedirs(save_path, exist_ok=True)
3     mem = set()
4     mem.add("/wiki/" + keyword)
5     dic_en = {}
6     for idx in tqdm(range(1, num + 1, 1)):
7         key = mem.pop()
8         try:
9             r = requests.get(url + key, headers=headers,
10                             proxies=proxies)
11             # r.encoding = r.apparent_encoding
12             soup = BeautifulSoup(r.text, 'html.parser')
13             text = ''
14             for x in soup.find_all('p'):
15                 text += x.text
16             cnt = 0
17             link = soup.find_all('a')
18             # random.shuffle(link)
19             for x in link:
20                 tmp = x.get('href')
21                 if tmp and tmp.startswith('/wiki/') and tmp not
22                 in mem:
23                     mem.add(tmp)
24                     cnt += 1
25                     if cnt == 3:
26                         break
27             dic_en[url + key] = text
28             r.close()
29         except:
30             pass
31         if idx % save_feq == 0:
32             u = list(dic_en.keys())
33             t = list(dic_en.values())
34             result = pd.DataFrame()
35             result["url"] = u
36             result["text"] = t
37             result.to_excel(os.path.join(save_path,
38                                         "./data{}.xlsx".format(idx)))
39             dic_en = {}
```

实验数据清洗

由于本实验只需要对中文汉字和英文字母进行相关计算，我们只需要使用正则表达式相关的库 `re` 过滤掉中文数据里的非汉字，统一英文数据的大小写，并且过滤到其中的非英文字母部分即可。最终清洗完毕，得到数据的样本规模约为： 7×10^7 个中文汉字， 2×10^8 个英文字母。

实验内容与结果分析

根据离散随机变量熵的公式 $H(x) = -\sum_{i=1}^n p(x_i) \log_2 p(x_i)$ ，我们利用 `collection` 库函数 `Counter` 统计中文汉字和英文字母出现的频率，构建一个键为汉字/字母，值为频率的字典，并用频率近似公式中的概率，便能够计算得到最终的熵。此外，为了分析数据规模对熵计算结果的影响，我们逐步扩大文本规模进行实验。具体代码实现如下：

```
1  # calculate entropy
2  def cal_h(dic, n):
3      h = 0
4      for v in dic.values():
5          h += - v / n * math.log(v / n, 2)
6      return h
7
8  # Chinese
9  for txt in tqdm(ch_txt_list):
10     if pd.isna(txt):
11         continue
12     txt_rm = re.sub(r'^[\u4e00-\u9fa5]+', '', txt)
13     ch_txt_list_processed.append(txt_rm)
14 ch_N = len(''.join(ch_txt_list_processed))
15
16 ch_scale = []
17 ch_H = []
18 for num in tqdm(range(350000, ch_N + 1, 350000)):
19     random.shuffle(ch_txt_list_processed)
20     txt_data = ''.join(ch_txt_list_processed)[:num]
21     d = dict(Counter(txt_data))
22     H_Chinese = cal_h(d, num)
23     ch_scale.append(num)
24     ch_H.append(H_Chinese)
25
26 # English
27 for txt in tqdm(en_txt_list):
28     if pd.isna(txt):
29         continue
30     try:
31         txt_rm = re.sub(r'^[\u0041-\u005a\u0061-\u007a]+', '',
txt)
```

```

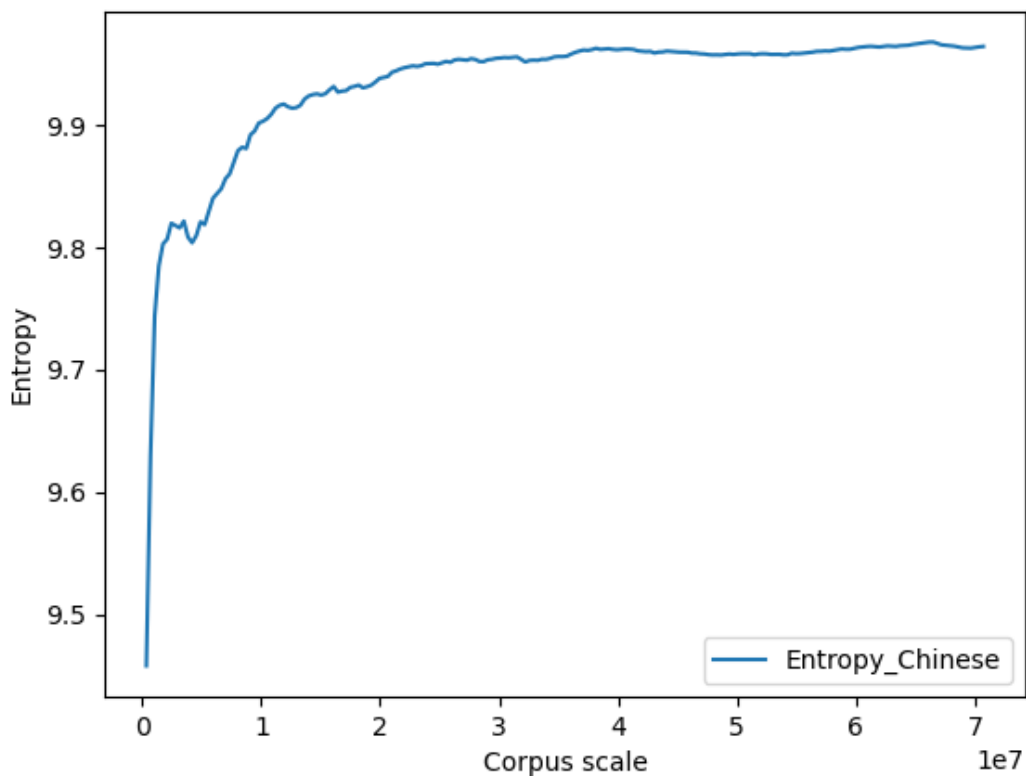
32     txt_rm = txt_rm.lower()
33     en_txt_list_processed.append(txt_rm)
34 except:
35     print(txt)
36 en_N = len("".join(en_txt_list_processed))
37
38 en_scale = []
39 en_H = []
40
41 for num in tqdm(range(1000000, en_N, 1000000)):
42     random.shuffle(en_txt_list_processed)
43     txt_data = "".join(en_txt_list_processed)[:num]
44     d = dict(Counter(txt_data))
45     H_English = cal_h(d, num)
46     en_scale.append(num)
47     en_H.append(H_English)

```

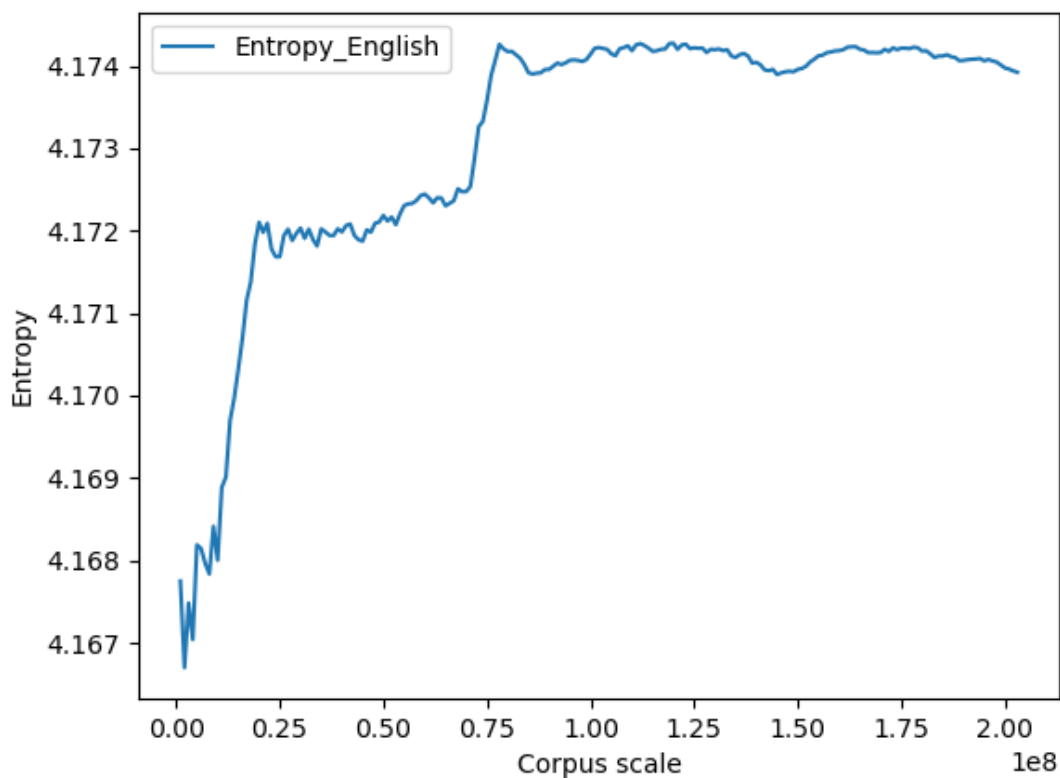
最终计算得到中文汉字和英文字母的熵如下：

$$H_{Chinese} = 9.96, H_{English} = 4.17$$

实验中，我们首先按收集到文本的顺序，每次增加总数据量的 $\frac{1}{20}$ ，得到如下的熵随文本数据规模的变化曲线：

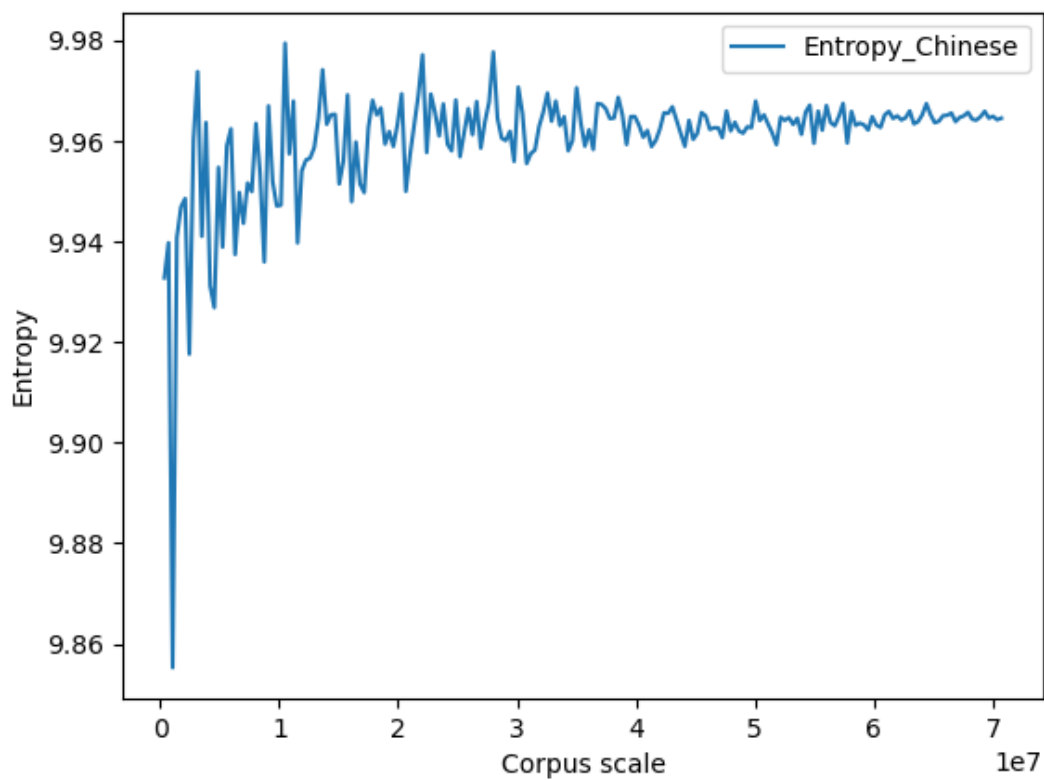


中文汉字的熵随文本数据规模扩大的变化曲线

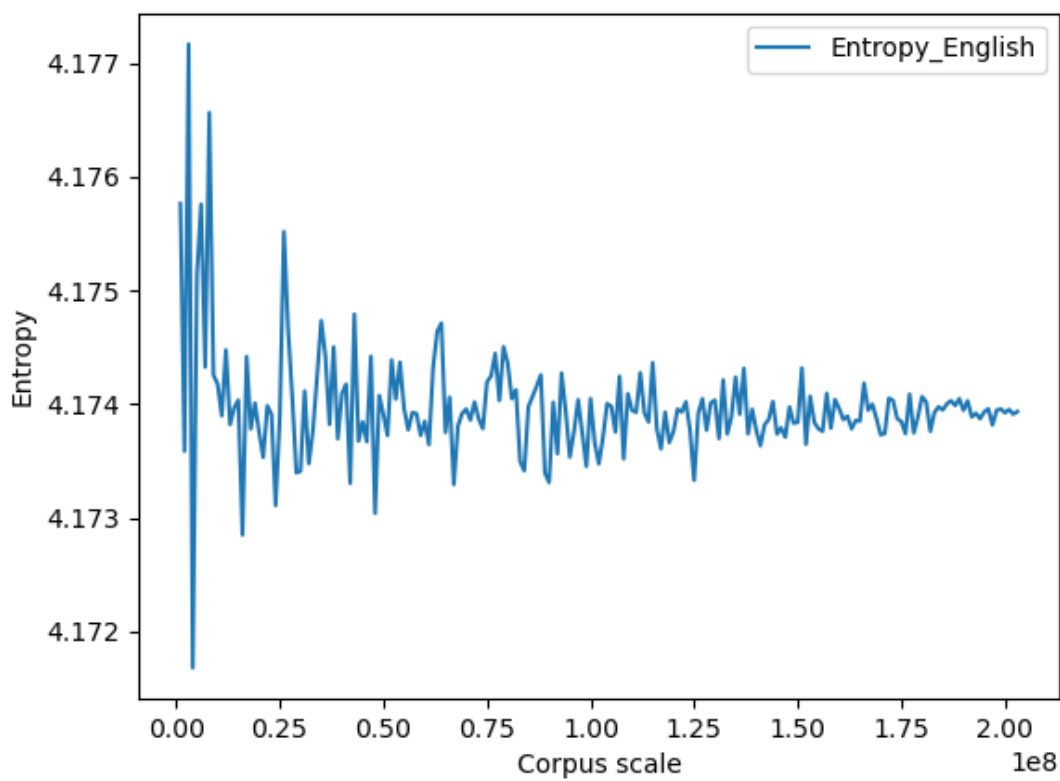


英文字母的熵随文本数据规模扩大的变化曲线

可以看出随着文本数据的增大，熵的值逐渐增加，最终呈现出收敛的趋势，这符合大数定理中频率趋近于概率的现象。此外，我们发现熵在数据规模刚开始扩大时有一个明显的上升趋势，而且整个过程波动较小。我们猜测这个现象与我们的数据收集方式有关。例如，我们的英文数据是通过 Wikipedia 词条之间的相互链接所爬取到的。因此，一开始收集到的数据之间相关性可能较大，从而熵较小。而随着数据规模的扩大，词条之间的相关性减弱，熵就呈现一个上升的趋势。我们的中文数据是通过 URL 链接中的 ID 每次加一爬取到的，相邻 ID 之前可能也有着较大的相关性。为了验证这一点，我们将原本的按照爬取顺序增加数据规模的方式改成随机打乱筛选的方式，得到如下结果：



中文汉字的熵随文本数据规模扩大的变化曲线



英文字母的熵随文本数据规模扩大的变化曲线

对比两次实验结果，我们发现在加入随机打乱筛选的操作后，熵的值波动增大，上升趋势减弱，并且依然呈现了收敛的趋势。