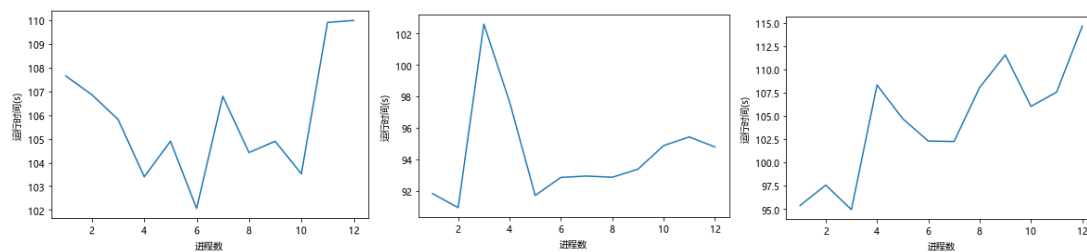


注：因 sogouCS 数据集过大，所以用某旅游评论数据集代替。

```
1  import time
2  import jieba
3  import pandas as pd
4  import matplotlib.pyplot as plt
5  from pathlib import Path
6  from multiprocessing import Process, Pool, Manager
7
8  plt.rcParams['font.sans-serif']=['Microsoft YaHei'] # 显示中文
9
10
11
12  def map_process(q, path):
13      # 读取文本数据
14      data = pd.read_csv(path).iloc[:, 0]
15      tokens = data.apply(lambda x: [i for i in jieba.lcut(x)], 1)
16      # 词频统计
17      all_words = (x for l in tokens for x in l)
18      count = pd.Series(all_words).value_counts()
19      q.put(count)
20
21  def reduce_process(q):
22      flag = 1
23      while True:
24          count = q.get()
25          if type(count).__name__ == 'NoneType':
26              break
27          else:
28              if flag:
29                  flag = 0
30                  C = count
31              else:
32                  for i in count.index:
33                      if i in C.index:
34                          C[i] += count[i]
35                      else:
36                          C[i] = count[i]
37      C.to_csv('word_count.csv')
38
39  if __name__ == '__main__':
40      run_time = []
41      for n in range(1, 13):
42          start = time.time() # 程序开始时间
43          q = Manager().Queue()
44          reduce = Process(target = reduce_process, args = (q,))
45          reduce.daemon = True
46          reduce.start()
47          maps = Pool(n)
48          for path in Path('data/').glob('*.txt'):
49              maps.apply_async(map_process, args=(q, path))
50          maps.close()
51          maps.join()
52          q.put(None) # 发出结束信号
53          reduce.join() # 等map结束后再结束reduce
54          end = time.time() # 程序结束时间
55          run_time.append(end - start)
56      plt.plot(list(range(1, 13)), run_time)
57      plt.xlabel('进程数')
58      plt.ylabel('运行时间(s)')
```

为消除运行结果的随机性，多次运行结果如下：



首先，可以看到每次运行的结果差异较大。设备 CPU 核数为 12，可能是电脑后台运行程序占用内存空间导致每次运行结果不一致，尤其是当进程数上升时运行时间反而上升。

总体上，当进程数为 1，4 或 6 时运行时间较短，效率较高，说明多进程的运行效率并不是进程数越多越好，受进程间通信时间的影响，运行效率有一个局部最优点。