

生信编程直播课程优秀学员作业展示2

Original x2yline 生信技能树 2017-03-17 12:06

收录于合集

#学徒作业

117个

题目：hg19基因组序列的一些探究

学员：x2yline

具体题目详情请参考生信技能树论坛

数据来源： <http://hgdownload.cse.ucsc.edu/goldenPath/hg19/bigZips/chromFa.tar.gz> 下载.gz数据后解压

R语言实现（太卡，高能报警）

代 码 地 址：
<https://raw.githubusercontent.com/x2yline/courseranotes/master/myscript/class2/fastafileGCandN.R>

代码内容：

```
1. setwd('E:\\r\\biotraineer_demo\\class 2')
2. # 读入数据
3. t1 <- Sys.time()
4. df <- read.csv('chr1.fa', header=F, stringsAsFactors=F)
5. # index_df 为chr所在的位置
6. index_df <- data.frame(begin=which(sapply(df[,1], function(x){
7.   substr(x, start=1, stop=1)=='>'})))
8. # index_df1 为string所在的位置+1
9. index_df1 <- data.frame(rbind(matrix(index_df[-1,1]),dim(df)[1]+1))
10. # 把index_start和index_end存入data.frame
11. index_df2 <- cbind(index_df, index_df1)
12. remove(index_df1, index_df)
13. # 得出每个染色体对应string后计算其N与GC百分比
14. result <- apply(index_df2, 1, function(x) { # 把提取字符串后把字符串变为大写
15.   y <- toupper(paste(df[(x[1]+1):(x[2]-1)],1, collapse=''))
16.   y <- strsplit(y, split=character(0))[[1]]
17.   N <- length(y[y=='N'])/length(y)
18.   GC <- length(y[y=='G' | y=='C'])/(length(y)-length(y[y=='N']))
19.   c(N,GC)
20. })
21. # 把行名改为N和GC并转秩
22. rownames(result) = c('N','GC')
23. result <- t(result)
24. # 取结果前几行
25. head(result)
```

```
26. difftime(Sys.time(), t1, units = 'secs')
```

由于电脑问题，试了一下1号染色体，电脑卡住了，于是又试了一下Y染色体，跑出来结果如下：

N	N ratio	GC	GC ratio	all base num
33720000	0.5679295	10252459	0.3996504	5933550

耗时：41.44945 secs

电脑配置信息：

- R version 3.3.2 (2016-10-31)
- Platform: x86_64-w64-mingw32/x64 (64-bit)
- Running under: Windows 7 x64 (build 7601) Service Pack 1

python3第一种实现方法（运行速度较快，但没有3快）

数据来源：<http://hgdownload.cse.ucsc.edu/goldenPath/hg19/bigZips/chromFa.tar.gz>

数据下载时间：2017-01-10 23:08

运行消耗时间：309 seconds

未优化速度的代码如下

```
1. import os
2. import time
3. begin = time.time()
4. os.chdir(r'F:\tmp\chromFa')
5. def count_n_and_gc(file):
6.     content = []
7.     chromosome = []
8.     g = 0; c = 0; n = 0; a = 0; t = 0
9.     with open(file) as f:
10.         raw_list = f.readlines()
11.         for i in raw_list:
12.             if not i.startswith('>'):
13.                 i = i.upper()
14.                 n += i.count('N')
15.                 g += i.count('G')
16.                 c += i.count('C')
17.                 a += i.count('A')
18.                 t += i.count('T')
19.             else:
20.                 if chromosome:
21.                     content.append((n, a, t, c, g))
22.                     g = 0; c = 0; n = 0; a = 0; t = 0
23.                     chromosome.append(i.strip())
24.                 content.append((n, a, t, c, g))
25.     return (content, chromosome)
```

```

26. content = []
27. chromosome = []
28. for i in (list(range(1,23)) + ['X','Y']):
29.     file = 'chr'+ str(i) + '.fa'
30.     print('Start dealing with ' + file)
31.     m, n = count_n_and_gc(file)
32.     content += m
33.     chromosome += n
34. all_info = 'chr,GC_ratio,N_ratio,Length,N,A,T,C,G'
35. for i in range(len(chromosome)):
36.     data = '\n'+str(chromosome[i]) + ',' + "%.5f"%((content[i][-1]+content[i][-2])/sum(content[i][1:])) + ',
37.     all_info += data
38. with open('hg19_analysis.csv','w') as f:
39.     f.write(all_info)
40. print('Time using:'+ str(time.time() - begin) + ' seconds\n')

```

shell + python3 (最快)

先使用shell脚本把所有chromFa.tar.gz 中的所有.fa文件合并为一个hg19.fa文件

脚本如下：

```

1. tar zvf chromFa.tar.gz
2. cat *.fa > hg19.fa
3. rm chr*.fa
4. less hg19.fa

```

按照老师的方法对python算法进行改良

改良后的代码如下：

代码地址：

```

1. import os
2. import time
3. import re
4. import sys
5. from collections import OrderedDict
6. start = time.clock()
7. def count_fasta_atcgn(file_path, buffer_size=1024*1024):
8.     bases = ['N', 'A', 'T', 'C', 'G']
9.     ATCG_analysis = OrderedDict()
10.    with open(file_path, 'r') as f:
11.        line1 = f.readline()
12.        chr_i = re.split('\s', line1)[0][1:]
13.        print(chr_i)
14.        ATCG_analysis[chr_i] = OrderedDict()
15.        for base in bases:
16.            ATCG_analysis[chr_i][base] = 0
17.        while True:
18.            chunk = f.read(buffer_size).upper()
19.            if '>' in chunk:
20.                chromosome = re.split('>', chunk)
21.                if chromosome[0]:
22.                    for base in bases:
23.                        ATCG_analysis[chr_i][base] += chromosome[0].count(base)
24.                for i in chromosome[1:]:
25.                    if i:
26.                        chr_i = re.split('\s', i[0:i.index('\n')])[0]
27.                        print(chr_i)
28.                        strings = i[i.index('\n')+1:].upper()

```

```

29.         ATCG_analysis[chr_i] = OrderedDict()
30.         for base in bases:
31.             ATCG_analysis[chr_i][base] = strings_i.count(base)
32.     else:
33.         for base in bases:
34.             ATCG_analysis[chr_i][base] += chunk.count(base)
35.     if not chunk:
36.         break
37.     return ATCG_analysis
38.
39. def write_atcg_to_csv(ATCG_analysis, file_path = '.'):
40.     file = os.path.join(file_path, 'atcg_analysis.csv')
41.     csv_content = 'chromosome\tGC_content\tN_content\tLength\tN\tA\tT\tC\tG\n'
42.     for chr_id, atcg_count in ATCG_analysis.items():
43.         GC = atcg_count['G'] + atcg_count['C']
44.         N = atcg_count['N']
45.         Length = sum(atcg_count.values())
46.         GC_content = GC*1.0/(Length-N)
47.         N_content = N*1.0/Length
48.         csv_content += chr_id + '\t' + '%.4f'%GC_content + '\t' + '%.4f'%N_content + '\t' + str(Length) +
49.     with open(file, 'w') as f:
50.         csv_file_content = re.sub('\t', ',', csv_content)
51.         f.write(csv_file_content)
52.     print(u'File have been saved in ' + file)
53.     return csv_content
54.
55. if sys.argv:
56.     result = OrderedDict()
57.     for f in sys.argv:
58.         done = 0
59.         f = f.strip('\'')
60.         if f.count('.') != 1 or f[-2:] == 'py' or not os.path.exists(f):
61.             continue
62.         print(f)
63.         try:
64.             done = 1
65.             result = OrderedDict(count_fasta_atcgn(file_path = f, buffer_size = 1024*2048), **result)
66.         except Exception as e:
67.             if f.startswith('-'):
68.                 pass
69.             else:
70.                 print(type(e))
71.         if done == 1:
72.             file = write_atcg_to_csv(result)
73.             print(file)
74.             print('used %.2f s'%(time.clock()-start))
75.         else:
76.             print('\n\nSorry! The command is invalid!\n')
77.     else:
78.         directory = input('Enter your file: ')
79.         start = time.clock()
80.         if directory.count('.') != 1 or directory[-2:] == 'py' or not os.path.exists(directory):
81.             print('Your file is invalid!')
82.         else:
83.             result = count_fasta_atcgn(file_path = directory, buffer_size = 1024*2048)
84.             file = write_atcg_to_csv(result)
85.             print('used %.2f s'%(time.clock()-start))

```

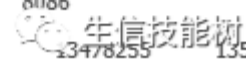
保存上述代码为 **fasta_atcgn_summary.py** 文件后

在命令行下输入：

```
1. python fasta_atcgn_summary.py F:\tmp\hg19.fa
```

部分输出结果如下

chromosome	GC_content	N_content	Length	N	A	T	C	G	
CHRUN_GL000237	0.4666	0	45867	0	12273	12191	10241	11162	
CHR14	0.4089	0.1776	107349540	19060000	25992966	26197495	18027132	1	
8071947									
CHR9_GL000199_RANDOM	0.3791	0	169874	0	54702	50765	34981	29426	
CHR2	0.4024	0.0205	243199373	4994855	71102632	71239379	47915465	479	
47042									
CHR8_GL000197_RANDOM	0.5401	0.0027	37175	100	8644	8408	9883	10140	
CHRUN_GL000247	0.436	0	36422	0	11002	9540	7794	8086	
CHR19	0.4836	0.0561	59128983	3320000	14390632	14428951	13478255	135	



使用python进一步进行可视化处理

代码如下:

```

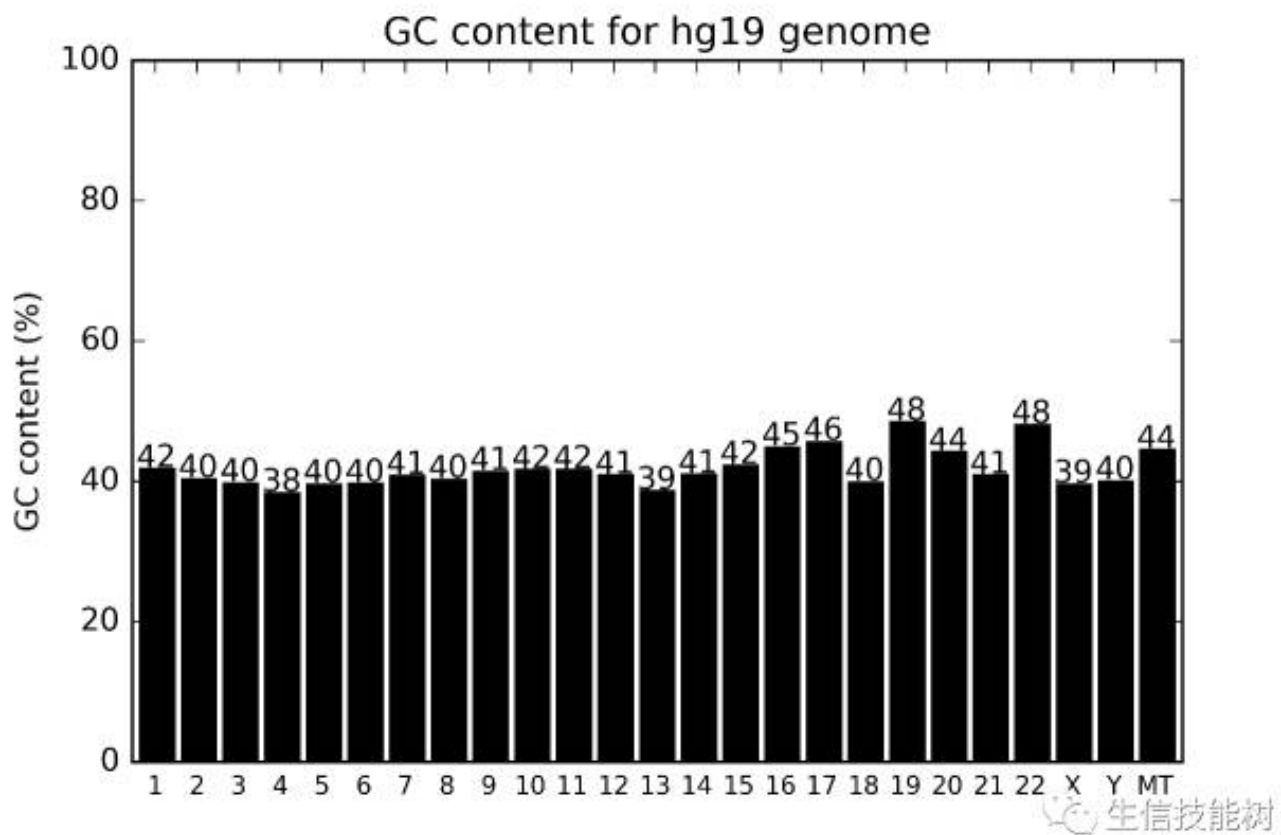
1. import os
2. import time
3. import re
4. import sys
5. from collections import OrderedDict
6. start = time.clock()
7. def count_fasta_atcgn(file_path, buffer_size=1024*1024):
8.     bases = ['N', 'A', 'T', 'C', 'G']
9.     ATCG_analysis = OrderedDict()
10.    with open(file_path, 'r') as f:
11.        line1 = f.readline().upper()
12.        chr_i = re.split('\s', line1)[0][1:]
13.        print(chr_i)
14.        ATCG_analysis[chr_i] = OrderedDict()
15.        for base in bases:
16.            ATCG_analysis[chr_i][base] = 0
17.        while True:
18.            chunk = f.read(buffer_size).upper()
19.            if '>' in chunk:
20.                chromosome = re.split('>', chunk)
21.                if chromosome[0]:
22.                    for base in bases:
23.                        ATCG_analysis[chr_i][base] += chromosome[0].count(base)
24.                for i in chromosome[1:]:
25.                    if i:
26.                        chr_i = re.split('\s', i[0:i.index('\n')])[0]
27.                        print(chr_i)
28.                        strings_i = i[i.index('\n'):]
29.                        ATCG_analysis[chr_i] = OrderedDict()
30.                        for base in bases:
31.                            ATCG_analysis[chr_i][base] = strings_i.count(base)
32.            else:
33.                for base in bases:
34.                    ATCG_analysis[chr_i][base] += chunk.count(base)
35.            if not chunk:
36.                break
37.    return ATCG_analysis
38.
39. def write_atcg_to_csv(ATCG_analysis, file_path = '.'):
40.    file = os.path.join(file_path, 'atcg_analysis.csv')
41.
42.    csv_content = 'chromosome\tGC_content\tN_content\tLength\tN\tA\tT\tC\tG\n'
43.    for chr_id, atcg_count in ATCG_analysis.items():
44.        GC = atcg_count['G'] + atcg_count['C']
45.        N = atcg_count['N']
46.        Length = sum(atcg_count.values())
47.        GC_content = GC*1.0/(Length-N)
48.        N_content = N*1.0/Length
49.        csv_content += chr_id + '\t' + '%.4f'%GC_content + '\t' + '%.4f'%N_content + '\t' + str(Length) + '\n'
50.    with open(file, 'w') as f:
51.        csv_file_content = re.sub('\t', ',', csv_content)
52.        f.write(csv_file_content)

```

```

52.     print(u'File have been saved in ' + file)
53.     return csv_content
54.
55. file_path = 'F:\genome\chromFa\hg19.fa'
56.
57. ATCG_analysis = count_fasta_atcgn(file_path, buffer_size=1024*1024)
58. cg_list = []
59. chr_id_list = list(range(1,23)) + ['X','Y','M']
60. for i in chr_id_list:
61.     cg_list.append((ATCG_analysis['CHR'+str(i)]['G']+ATCG_analysis['CHR'+str(i)]['C'])/(ATCG_analysis['CHR'+str(i)]['G']+ATCG_analysis['CHR'+str(i)]['C']+ATCG_analysis['CHR'+str(i)]['A']+ATCG_analysis['CHR'+str(i)]['T']))
62. import matplotlib.pyplot as plt
63. plt.bar(left = range(25), height = cg_list, color='k')
64. for i in range(len(cg_list)):
65.     plt.text( x=i-0.1, y=cg_list[i]+.35,s=str(round(cg_list[i])))
66. plt.title('GC content for hg19 genome')
67. plt.ylabel('GC content (%)')
68. pos = []
69. for i in range(len(chr_id_list)):
70.     pos.append(i + 0.35)
71. plt.xticks(pos, list(range(1,23)) + ['X','Y','MT'], fontsize=8)
72. plt.xlim(-0.2, )
73. plt.ylim(0, 100)
74. plt.savefig('F:\hg19_gc.png',dpi=600)
75. plt.show()

```



本文编辑：思考问题的熊



收录于合集 #学徒作业 117

上一篇

给学徒的GEO作业

下一篇

生信编程直播课程优秀学员学习心得及作业展示3

People who liked this content also liked

EDCI, HATU, CDI等酸胺缩合试剂副产物和中间态结构和MS分析
有机合成路线



TLC点板，产物极性忽大忽小，咋回事！
有机合成路线



人肺内皮细胞整合单细胞图谱
生信菜鸟团

