

# 第四次上机实验

生信 2001 张子栋 2020317210101

GitHub 地址: [MarkdownNotes/R at main · Bluuur/MarkdownNotes \(github.com\)](https://github.com/Bluuur/MarkdownNotes)

- 本次实验所需文件
  - hg19\_gene\_table.txt
  - hg19.chrom.sizes.txt

要求: 不用任何循环语句

1. 对人染色体长度数据 hg19.chrom.sizes.txt

1. 用 read.table 正确读入数据

```
1 data <- read.table('hg19.chrom.sizes.txt', header = F)
```

2. 最长和最短的染色体分别是哪条? 各多长?

```
1 data <- read.table('hg19.chrom.sizes.txt', header = F)
2 data[which(data$V2 == max(data$V2)),]
3 data[which(data$V2 == min(data$V2)),]
```

A data.frame: 1 × 2

	V1	V2
	<chr>	<int>
1	chr1	249250621

A data.frame: 1 × 2

	V1	V2
	<chr>	<int>
22	chr21	48129895

3. 求所有染色体的总长度和平均长度

```
1 data <- read.table('hg19.chrom.sizes.txt', header = F)
2 cat('Length in total', sum(data$V2), '\n')
3 cat('Average length', mean(data$V2))
```

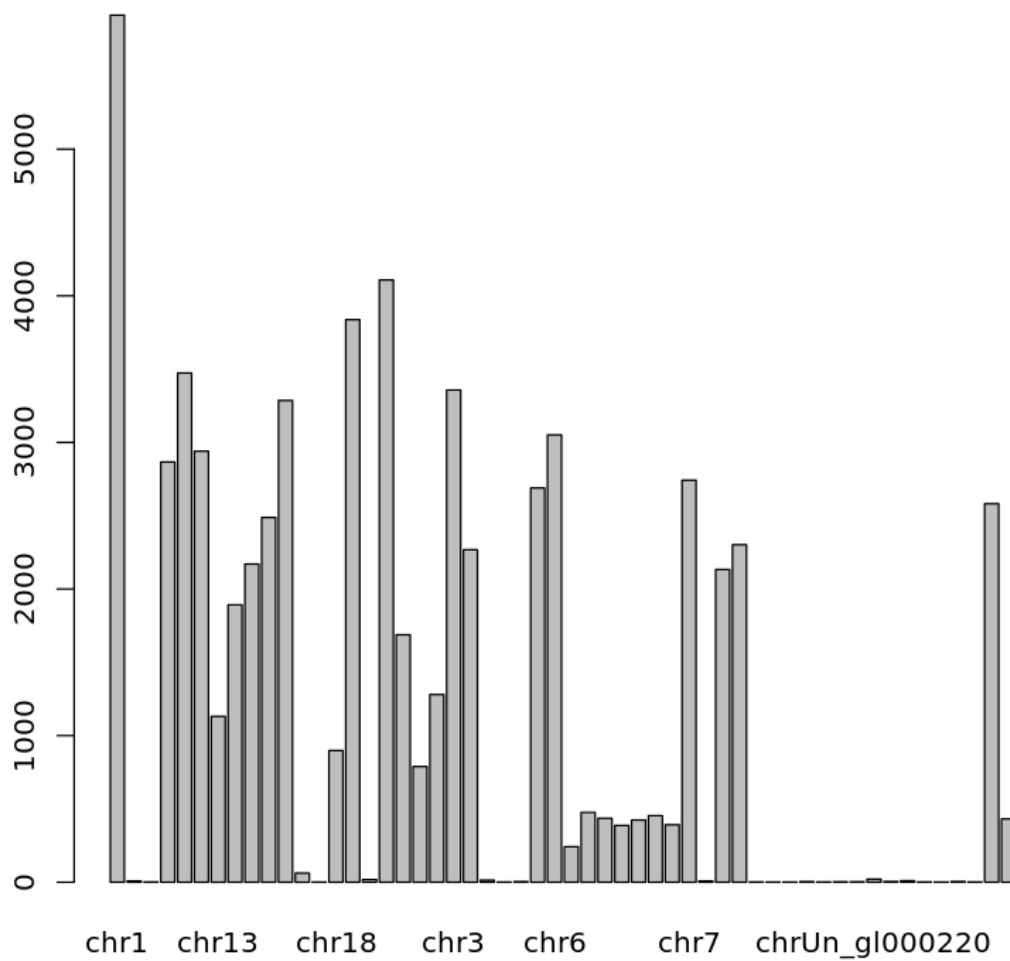
```
1 Length in total 3095677412
2 Average length 128986559
```

2.用 `read.table` 读入 `hg19.chrom.sizes.txt`

1.用函数 `table` 得到每条染色体上的基因个数, 并画 `barplot`

```
1 data <- read.table("hg19_gene_table.txt", header = T)
2 table(data$chrom)
3 barplot(table(data$chrom))
```

1		chr1	chr1_g1000191_random	chr1_g1000192_random
2		5914	9	2
3		chr10	chr11	chr12
4		2866	3473	2940
5		chr13	chr14	chr15
6		1131	1892	2170
7		chr16	chr17	chr17_ctg5_hap1
8		2488	3286	62
9	chr17_g1000205_random		chr18	chr19
10		1	898	3837
11	chr19_g1000209_random		chr2	chr20
12		18	4108	1688
13		chr21	chr22	chr3
14		789	1280	3357
15		chr4	chr4_ctg9_hap1	chr4_g1000193_random
16		2268	14	1
17	chr4_g1000194_random		chr5	chr6
18		4	2689	3051
19	chr6_apd_hap1	chr6_cox_hap2	chr6_dbb_hap3	
20		242	476	436
21	chr6_mann_hap4	chr6_mcf_hap5	chr6_qb1_hap6	
22		387	424	454
23	chr6_ssto_hap7	chr7	chr7_g1000195_random	
24		392	2742	9
25		chr8	chr9	chrM
26		2133	2302	2
27	chrUn_g1000211	chrUn_g1000212	chrUn_g1000213	
28		1	1	4
29	chrUn_g1000215	chrUn_g1000218	chrUn_g1000219	
30		2	3	3
31	chrUn_g1000220	chrUn_g1000222	chrUn_g1000223	
32		21	5	10
33	chrUn_g1000224	chrUn_g1000227	chrUn_g1000228	
34		2	1	5
35	chrUn_g1000241	chrX	chrY	
36		2	2582	432



2.结合染色体长度计算每条染色体上的基因密度(每 Mbp 的基因个数)

```

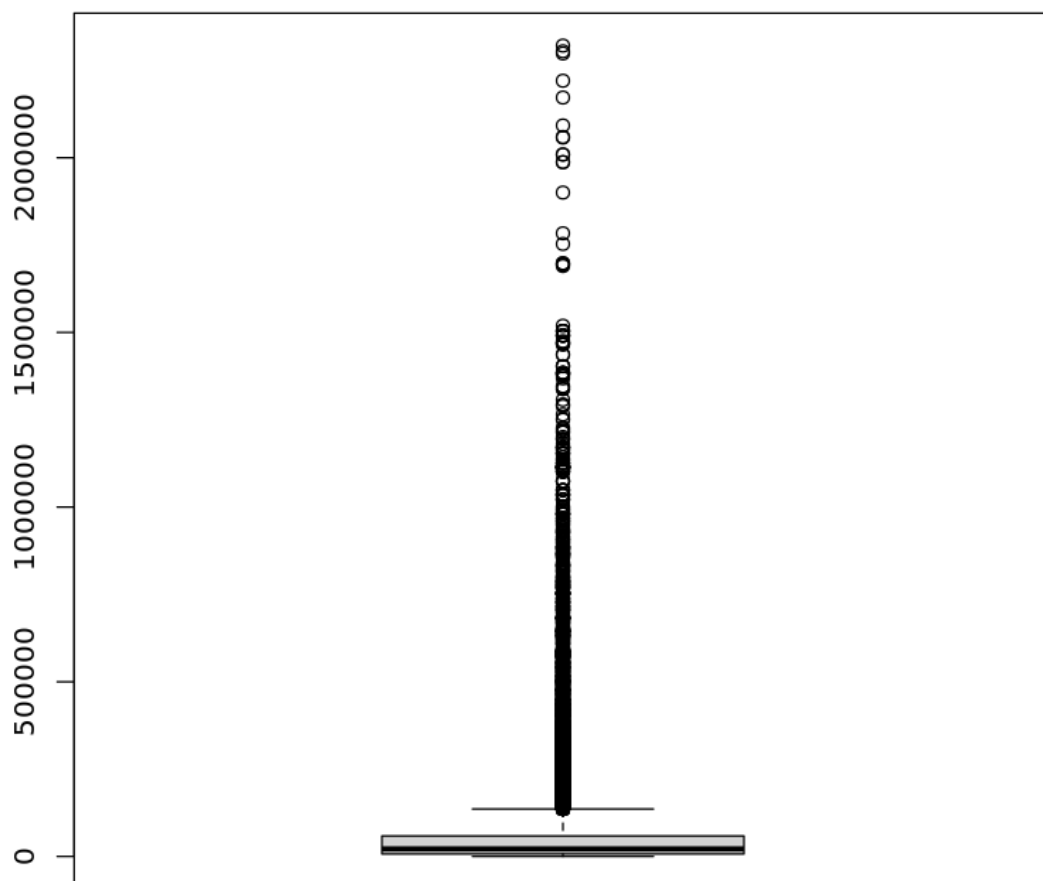
1 data <- read.table("hg19_gene_table.txt", header = T)
2 chromoLength <- read.table('hg19.chrom.sizes.txt', header = F)
3 geneNum <- table(data$chrom)
4 geneNum <- geneNum[as.vector(chromoLength$V1)]
5 geneNum / (chromoLength$V2 / 1000000)

```

1	chr1	chr2	chr3	chr4	chr5	chr6	chr7
	chr8						
2	23.727122	16.891491	16.952625	11.864762	14.863312	17.830107	17.230257
	14.573254						
3	chr9	chr10	chr11	chr12	chr13	chr14	chr15
	chr16						
4	16.301566	21.145869	25.724684	21.964575	9.820276	17.624668	21.164250
	27.535906						
5	chr17	chr18	chr19	chr20	chr22	chr21	chrX
	chrY						
6	40.470368	11.501430	64.892034	26.782802	24.949046	16.393138	16.629038
	7.275965						

3.计算基因长度, 绘制基因长度的 `boxplot`

```
1 data <- read.table("hg19_gene_table.txt", header = T)
2 geneLength <- data$txEnd - data$txStart + 1
3 boxplot(geneLength)
```



4.找出最长的基因

```

1 data <- read.table("hg19_gene_table.txt", header = T)
2 gene <- cbind(data$geneName, data$txEnd - data$txStart + 1)
3 gene[which(gene[,2]==max(gene[,2])),]

```

'ATL1' · '99986'

5. 阅读 `cor` 函数的帮助文件, 计算基因长度和外显子个数的相关系数

```

1 data <- read.table("hg19_gene_table.txt", header = T)
2 gene <- cbind(data$txEnd - data$txStart + 1, data$exonCount)
3 cor(gene[,1], gene[,2])

```

0.376018862127503

6. 用 `prop.table` 分染色体计算 +/- strand 上基因的百分比

```

1 data <- read.table("hg19_gene_table.txt", header = T)
2 gene <- table(data$strand)
3 prop.table(gene)

```

```

1      -      +
2 0.4868818 0.5131182

```