利用 Python 进行生信分析 | 基础 + 实战

——Python for Bioinformatics



赵华男

2022 一起加油

项目一 序列文件的处理

项目一序列文件的处理

序列储存格式的介绍

序列储存格式的介绍-FASTA 文件格式

• FASTA 格式

- 一种用于表示核苷酸序列或多肽序列的文本格式;其中碱基对或氨基酸用单个字母来表示
- 允许在序列前添加序列名及注释
- 该格式已成为生物信息学领域的一项标准
- FASTA 文件各行记录信息如下:
 - 第一行
 - 由大于号">" 开头的任意文字说明,用于序列标记
 - 为了保证后续分析软件能够区分每条序列,单个序列的标识必须是唯一的
 - 第二行
 - 序列本身: 只允许使用既定的核苷酸或氨基酸编码符号。
 - 通常核苷酸符号大小写均可, 而氨基酸常用大写字母。
 - 注意有些程序对大小写有明确要求。一般每行 60-80 个字 母.

>HWI-D00433:463:HNT7JBCXX:1:1101:19071:2193 1:N:0:TTCTCCAT

CCTACGGGGTTCACCAGTAGGGATCTTCCACAATGGGCGAAAGCCTGATTGAGCAAACGCGGTTTTTAAGAAGGTCTTCGGATCGTAAAACCCTGTTGTTAGAGAAGAAAGTGCGCGGCTAACTACTTCTACTGTACTAACAAGAAAGCACCGGCTAACTACGTTC

序列储存格式的介绍-FASTA 文件格式

```
# legend server
    cd /home/zhaohuanan/3.project/2022 Other projects/2022-09-25 Prepared data/FASTA
    rsync -avzP rsync://hgdownload.cse.ucsc.edu/goldenPath/mm39/chromosomes/ .
    # check - T md5
    md5sum -c md5sum txt
 6
7
    (base) PS C:\Users\vagrant\PythonForBioinformatics> rsync -avzP rsync://hgdownload.cse.
    ucsc.edu/goldenPath/mm39/chromosomes/chrM.fa.gz .
8
    receiving incremental file list
10
    chrM.fa.gz
              5.385 100%
11
                          5.14MB/s
                                       0:00:00 (xfr#1, to-chk=0/1)
12
13
    sent 43 bytes received 5,492 bytes 299.19 bytes/sec
    total size is 5.385 speedup is 0.97
14
15
    #接下来我们需要将genome生成一下,未确定的基因组区域暂不考虑,只看确定的基因组
16
17
    # 测试命今
    echo `seg 1 1 19` X Y M | awk 'BEGIN {printf "cat "}{for(i=1: i<=NF:i++){printf "chr"$i".fa.gz "}}'
18
19
    # 写入 genome 文件
20
    echo `seg 1 1 19` X Y M | awk 'BEGIN {printf "cat "}{for(i=1: i<=NF:i++){printf "chr"$i".fa.gz "}}' | \
        sh > genome ucsc mm39.fa.gz
21
    zcat genome_ucsc_mm39.fa.gz | grep -v N | less # linux/window
22
23
    zcat < genome ucsc mm39.fa.gz | grep -v N | less # macos
```

序列储存格式的介绍-FASTQ 文件格式

第一行

- 以 @ 开头
- 后面是 reads 的 ID 以及其他信息
- 例如下一页例中 HWUSI-EAS100R 代表 Illumina 设备名称,
- 6 代表 flowcell 中的第六个 lane.73 代表第六个 lane 中的第 73 个 tile
- 941:1973 代表该 read 在该 tile 中的 x:v 坐标信息:
- #0. 若为多样本的混合作为输入样本. 则该标志代表样本的编号. 用来区分个样本中的 reads/
- /1 代表 paired end 中的前一个 read。

• 第二行

- read 的序列
- 紧接着下面两行代表该 read 的质量

• 第三行

● 以"+"开头, 跟随着该 read 的名称 (一般于 @ 后面的内容相同), 但有时可以省略, 但"+"一定不能省

https://baike.baidu.com/item/fastQ 格式 (rmspace)

序列储存格式的介绍-FASTQ 文件格式

第四行

- 代表 reads 的质量。这一行可以详细说一下!
- Illumina 测序仪是按照荧光信号来判断所测序的碱基是哪一种的,例如红黄蓝绿分别对应 ATCG,那么一旦出现一个紫色的信号该怎么判断呢?
- 因此对每个结果都有一个概率的问题。起初 sanger 中心用 Phred quality score 来衡量该 read 中每个碱基的质量, 既-10lgP #(lg 意为 log10)
 - 其中 P 代表该碱基被测序错误的概率
 - 如果该碱基测序出错的概率为 0.001, 则 Q 应该为 30, 那么 30+33=63, 那么 63 对应的 ASCii 码为 '?', 则在第四行中该碱基对应的 质量代表值即为 '?'
 - ASCII https://baike.baidu.com/item/ASCII

序列储存格式的介绍-FASTQ 文件格式

@E00591:528:HHVW3CCX2:2:1101:17360:2170 1:N:0:GATCAGCG:ACTA

from illumina

15

FASTQ 文件的操作

FASTQ 文件的操作

● FASTQ 文件的操作

- 读取 FASTQ 文件并以 FASTA 格式输出
- 解析 FASTQ 的质量值, 计算 Q30 比例
- 根据 FastQC 报告对 FASTQ 文件进行截取
- 根据 FastQC 报告, 过滤低质量的 lane,tile 数据

FASTA 文件的操作

FASTA 文件的操作

● FASTA 文件的操作

- 读取 FASTA 文件、并将其中 U 替换成 T
- 读取 FASTA 文件,并输出反向互补序列
- 计算基因组序列的长度
- 计算基因组各染色体的平均 GC 含量
- 计算基因组中 N 的总长度 (effective length)

Docstring-reST

```
def func(arg1, arg2):
         """Summary line.
         Extended description of function.
         :param int argl: Description of argl.
         :param str arg2: Description of arg2.
         :raise: ValueError if arg1 is equal to arg2
 8
         :return: Description of return value
 9
10
         :rtvpe: bool
11
12
         :example:
13
14
         >>> a=1
        >>> h=2
15
16
         >>> func(a,b)
17
         True
         ....
18
19
20
         if arg1 == arg2:
             raise ValueError('arg1 must not be equal to arg2')
21
22
23
         return True
```

Docstring-Google Style

```
def func(arg1, arg2):
    """Summary line.
    Extended description of function.
    Aras:
        arg1 (int): Description of arg1
        arg2 (str): Description of arg2
    Returns:
                                                                          if arg1 == arg2:
        bool: Description of return value
                                                                              raise ValueError(
                                                                  3
                                                                                  'arg1 must not be equal to arg2'
    Raises:
        AttributeError: The ''Raises'' section is a list of all
            exceptions that are relevant to the interface.
                                                                  6
                                                                          return True
        ValueError: If 'arg2' is equal to 'arg1'.
    Examples:
        Examples should be written in doctest format, and should
            illustrate how to use the function.
        >>> a=1
        >>> h=2
        >>> func(a,b)
        Truo
```

10

11

12

13

14 15

16

17 18

19

20

21

22

23

24

25

....

Docstring-Numpy Style

```
. . . .
def func(arg1, arg2):
                                                                            ValueError
    """Summary line.
                                                                                If 'arg2' is equal to 'arg1'.
    Extended description of function.
                                                                            See Also
    Parameters
                                                                            otherfunc: some other related function
    arg1 : int
                                                                            Examples
        Description of arg1
                                                                   10
    arg2 : str
                                                                   11
                                                                            These are written in doctest format, and should
        Description of arg2
                                                                   12
                                                                            use the function.
                                                                   13
    Returns
                                                                   14
                                                                            >>> a=1
                                                                   15
                                                                            >>> h=2
                                                                            >>> func(a,b)
                                                                   16
        Description of return value
                                                                   17
                                                                            True
                                                                             ....
                                                                   18
    Raises
                                                                   19
                                                                   20
                                                                            if arg1 == arg2:
    AttributeError
                                                                                raise ValueError(
                                                                   21
        The ''Raises'' section is a list of all exceptions
                                                                   22
                                                                                     'arg1 must not be equal to arg2'
        that are relevant to the interface.
                                                                   23
    ....
                                                                   24
                                                                            return True
```

10

11

12

13

14

15

16

17

18

19

20

21

22

23

项目实战一 FASTA 与 FASTQ

STQ FASTA 文件的操作

Q & A!

Thank you!