

fastx-utils之去除重复序列: uniques

一、fastx-utils uniques介绍

功能描述:

`fastx-utils uniques` 根据序列去除重复。

命令行接口:

```
1 $ fastx-utils uniques
2
3 Usage: fastx-utils uniques [options] <in.fa/q>
4
5 Options:
6   -l STR  sequence label, default: [Uniq]
7   -m      print membership to stderr.
```

可选参数:

```
1  -l      序列标签, 默认为 'Uniq';
2  -m      输出标准错误;
```

二、使用场景实例及其用法

使用场景经典案例:

1. 16S扩增子数据分析: 对合并后的PE数据进行去除重复序列, 并进行计数;

示例演示:

示例文件: `A1.fastq`

```
1 $ zcat A1.fastq.gz | head -n 8
```

```

1 @HISEQ:483:HLJ2LBCXY:1:1101:7924:2136
2 TGGGGAATATTGGACAATGGGCGCAAGCCTGATCCAGCCATGCCGCGTGAGTGATGAAGGCCCTAGGGTTGTAAAGC
  CCTTTTCGGCGGGGAAGATAATGACGGTACCCGCGAGAAGAAGCCCCGGCTAACTTCGTGCCAGCAGCCGCGGTAATAC
  GAAGGGGGCTAGCGTTGCTCGGAATTACTGGGCGTAAAGCGCACGTAGGCGGCTTTCTAAGTCGGGGGTGAACTCCT
  GGAGCTTAACTCCAGAACTGCCTTCGATACTGGAGAGCTCGAGTCCGGGAGAGGTGAGTGGAAGTGCAGGTGTAGAG
  GTGAAACTCGTAGATATTTCGAAGAACACCAGTGGCGAAGGCGGCTCACTGGCCCCGCTACTGACGCTGAGGTGCGAA
  AGCGTGGGGAGCAAACAGG
3 +
4 <DEG/CEHHIHEHCHECFHICHHHHDHIICHHIEHCHCEGHHIICH?EEHFHHHIECEE?HE?
  1GHHHHHDHHHHFFGHHHCHHHII?H?HHEHCHF1CGHC0DEHEHHD<<FEHHHGHIC??
  HHCGHHHHIHH@FFEEHIDHHCHHHIHIHE@GH-G??EDEGEHIIHICHHC?EHH-
  KKKGK=KKKKKKKKKKKKKKKKKKHE@@G@CEC6+>=>8.-C@B8-HHHHACHCBFA.GA?
  HHEHHFHHCIEH@CGBAFGIHHHHCEHHCC,HEHHHGDHFFHEGCHEEHH@HFGHGIHHHHHECEHGIIEH
  HCEC/C/DGIHHGCC@HGD00C01GCHCCDE?C1D1GCC?DH=HFHHG@CEDEGIHHIFE=EHEHHHCHHHHHG?
  EHEG@<1
5 @HISEQ:483:HLJ2LBCXY:1:2108:14400:12517
6 TGGGGAATATTGGACAATGGGCGCAAGCCTGATCCAGCCATGCCGCGTGGGTGATGAAGGCCCTAGGGTTGTAAAGC
  CGTTTCGGCGGGGAAGATAATGACGGTACCCGCGAGAAGAAGCCCCGGCTAACTTCGTGCCAGCAGCCGCGGTAATAC
  GAAGGGGGCTAGCGTTGCTCGGAATTACGGGGCGTAAAGCGCACGTAGGCGACTCTTTAAGTCGGGGGTGAAATCCT
  GGAGCTCAAACCCAGAACTGCCTTCGATACTGGGAGAGCTCGAGTTCGGGAGAGGTGAGTGGAAGTGCAGGTGTAGCG
  GTGAAATTCGTAGATATTTCGAAGAACACCAGTGGCGAAGGCGGCTCACTGGCCCCGATACTGACGCTGAGGTGCGAA
  AGCGTGGGCAGCAAACAGG
7 +
8 HHHIIGHIIIIIIIGIIIIHHGIIIIIIHIGIIHIGIIIIIIIIIIIIHICHHHIIIIIIIIHIIHHHHIIIIIIH
  G/<D@GDHIIHH?HEFHHHHIEFEHEHHHHHDHHIIIIHH?BAHHDHHHFHHIIIIIGHII@@FHHHDHHH?
  AHHA?GHHDHGIIHHIIHHIIII<HHIIHI.87<DHHID?HHKKKKKKKKKK$K<KKKKKKKKKKKKKKKK-
  4AA.B8.88.-88..8.HHAB6.EHCHGCB@@A=HHHG.?GCF.</DCEDGEH?
  IIFHFG@HHHEHHHDHFCGCE.CC0IIIIHHCGHHEHHHHHDCHGEIHHIIHHHHIIHECHH@HHDHH@F?
  IHHHHHHHHIIHHHHIHHHF@C<<0<IIHHHHEIHD1HHHHHHIHHF

```

运行命令：去除 A1.fastq 文件中的重复序列, 使用 `-l Uniq` 重新命令.

```

1 $ fastx-utils uniques -l Uniq A1.fastq.gz | head -n 6
2 >Uniq_1
3 TGGGGAATATTGGACAATGGGCGCAAGCCTGATCCAGCCATGCCGCGTGAGTGATGAAGGCCCTAGGGTTGTAAAGC
  CCTTTTCGGCGGGGAAGATAATGACGGTACCCGCGAGAAGAAGCCCCGGCTAACTTCGTGCCAGCAGCCGCGGTAATAC
  GAAGGGGGCTAGCGTTGCTCGGAATTACTGGGCGTAAAGCGCACGTAGGCGGCTTTCTAAGTCGGGGGTGAACTCCT
  GGAGCTTAACTCCAGAACTGCCTTCGATACTGGAGAGCTCGAGTCCGGGAGAGGTGAGTGGAAGTGCAGGTGTAGAG
  GTGAAACTCGTAGATATTTCGAAGAACACCAGTGGCGAAGGCGGCTCACTGGCCCCGCTACTGACGCTGAGGTGCGAA
  AGCGTGGGGAGCAAACAGG
4 >Uniq_2
5 TGGGGAATATTGGACAATGGGCGCAAGCCTGATCCAGCCATGCCGCGTGGGTGATGAAGGCCCTAGGGTTGTAAAGC
  CGTTTCGGCGGGGAAGATAATGACGGTACCCGCGAGAAGAAGCCCCGGCTAACTTCGTGCCAGCAGCCGCGGTAATAC
  GAAGGGGGCTAGCGTTGCTCGGAATTACGGGGCGTAAAGCGCACGTAGGCGACTCTTTAAGTCGGGGGTGAAATCCT
  GGAGCTCAAACCCAGAACTGCCTTCGATACTGGGAGAGCTCGAGTTCGGGAGAGGTGAGTGGAAGTGCAGGTGTAGCG
  GTGAAATTCGTAGATATTTCGAAGAACACCAGTGGCGAAGGCGGCTCACTGGCCCCGATACTGACGCTGAGGTGCGAA
  AGCGTGGGCAGCAAACAGG
6 >Uniq_3
7 TGGGGAATATTGGACAATGGGCGAAAGCCTGATCCAGCCATGCCGCGTGGATGAAGGAGGCCCTAGGGTTGTAAAGT
  CCTTTTCGATGGTGAAGATAATGACGGTAACCATACAAGAAGCCCCGGCTAATTTTCGTGCCAGCAGCCGCGGTAATAC
  GAAAGGGGCTAGCGTTGCTCGGAATTACTGGGCGTAAAGCGCACGTAGGCGGCGGTGTAAGTTGGGGGTGAGATCCC
  GGGGCTCAACCCGGAAGTGCCTCCAAAACATAGCTAGAGGATGTGAGAGGACAGTGGAATTCGAGTGTAGAG
  GTGAAATTCGTAGATATTTCGAAGAACACCAGTGGCGAAGGCGACTGTCTGGCACATTTCTGACGCTGAGGTGCAA
  AGCGTGGGGAGCAAACAGG

```

本文材料为 **BASE (Biostack Applied bioinformatic SEies)** 课程 **Linux Command Line Tools for Life Scientists** 材料， 版权归 上海逻捷信息科技有限公司 所有。

Last Update: 2020-08-10 11:56 AM