

fastx-utils之根据标识符去除重复序列：dedup

一、fastx-utils dedup 介绍

功能描述：

`fastx-utils dedup` 根据序列去除重复。

命令行接口：

```
1 $ fastx-utils dedup
2
3 Usage: fastx-utils dedup <in.fa/fq>
```

二、使用场景实例及其用法

示例演示：

示例文件： `protein.faa`

```
1 $ cat protein.faa | head -n 10
```

```
1 >AGE80385.1
2 MKKIIIIISATTIVIGITSFAYFGSKTPLHNEAKAVESQKHNNHKKEEIPAFPKADHNAKKIDNDFS SVVTNP KSNLV
  LINKHRKLPDGYPEDLTRPNVPFISPKDKEKTL LRKDAAEALENMFKAAKKEGLDLTAVSGYRSYKRQKSLHDTY
  VRRQ GKAEANSVSAIPGTSEHQTGLAMD ISSSAKFQLEPIFGETAEGKWVAEHAHEFGFVIRYLEDKTDTTEYAY
  EPWHLRYVGNPYATYLYKHHLTLEEAMEDKK
3 >AGE79558.1
4 MINHELVERINFLAKKAKAEGLTEEEQRERQSLREQYLKGFRQNMLNELKGIKVVNEEGTDVTP TKLKALKKKQDNA
  KLN
5 >AGE81073.1
6 MDKVISNEILQQFKDRMRLGDDEDANLRRILFASNKDLTRVCGNYDLNIDEVFKELVFERSRYVYNDALEYFDKNF
  LSQINLSIGKALEAIKLDGD
7 >AGE80385.1
8 MQNGKVKWFNSEKGFGFIEVEGGEDVFVHFSAIQGEGFKTLEEQEVTFEVEQGNRGPQATNVNKK
9 >AGE81522.1
10 MKKTILTTTAALTMMGTGMGINVDHIKPAEVKADTISFYDVPNNHWATKAITNLANRNIVVGYGNGQFGFGDNVTR
    GQVARMIIYNLKPADAGNFKNPFSDIKGHMFEKEILALAKVGIIKGYGEGKFGPDDILTREQMAQVLTNAFKFEGT
    KKTSFVDVDKNSWSYKAIGALEEKGV TIGTGGNMYSPTS SVTREQYSQFLFNSINVIEKETKPEEK PNTGGGEVKPE
    EKPNTGGGEVKPEEK PNTGEETKPVNIPEWLETSLATNDFTFQAWYDGSEAINKAASTNAQQIVKNINSKYGTNLK
    YSEVGAI VQLVDGAREQLWLAGMNVNDFRVTFRVSNNAMIELTKELVTLVNSDLNLDQEIQEIP SAPMKIKNVEKG
    DYKIRISPAMADQMIIIIIEKK
```

运行命令： 去除 `protein.faa` 文件中的重复序列, 使用 `-l uniq` 重新命令.

```
1 $ fastx-utils dedup protein.faa | head -n 10
```

```
1 >AGE80385.1
2 MKKIIIIISATTIVIGITSFAYFGSKTPLHNEAKAVESQKHNNHKKEEIPAFPKADHNAKKIDNDFS SVVTNPKSNLV
  LINKHRKLPDGYIPEDLTRPNVPFISPKDKEKTLRKDAAEALENMFKA AKKEGLDLTAVSGYRSYKRQKSLHDTY
  VRRQGKAEANSVSAIPGTSEHQTGLAMDISSKSAKFQLEPIFGETAEGKWVAEHAHEFGFVIRYLEDKTDTTEYAY
  EPWHLRYVGNPYATYLYKHHLTLEEAMEDKK
3 >AGE79558.1
4 MINHELVERINFLAKKAKAEGLTEEEQRERQSLREQYLKGFRQNMLNELKGIKVVNEEGTDVTPTKLKALKKQDNA
  KLN
5 >AGE81073.1
6 MDKVISNEILQQFKDRMRLGDDEDANLRRILFASNKDLTRVCGNYDLNIDEVFKELVFERSRYVYNDALEYFDKNF
  LSQINLSIGKALEAIKLDGD
7 >AGE81522.1
8 MKKTILTTTAALTMMGTGMGINVDHIKPAEVKADTISFYDVPNNHWATKAITNLANRNIVVGYGNGQFGFGDNVTR
  GQVARMIIYNLKPADAGNFKNPFSDIKGHMFEKEILALAKVGIKGYGEGKFGPDDILTREQMAQVLTNAFKFEGT
  KKTSFVDVDKNSWSYKAIGALEEKGVITIGTGGNMYSPTSVVTREQYSQFLFNSINVIEKETKPEEKPN TGGEVKPE
  EKPNTGGEVKPEEKPN TGEETKPVNIPEWLETSLATNDFTFTQAWYDGSEAINKAASTNAQQIVKNINSKYGTNLK
  YSEVGAIQQLVDGAREQLWLAGMNVNDFRVTFRVSNNAMIELTKELVTLVNSDLNLDQEIQEIPSAPMKIKNVEKG
  DYKIRISPAMADQMITIIIEKK
9 >AGE80137.1
10 MKLLDLLSKGIVIGDGAVGTLLHSHGLQSSFEELNVSDPDLIISIHKQYVAAGADVIQTNTYGANEA KL RMYGLEN
  QVTKINRAAVKLAKASVTDKNAILGTIGGMKHIGAVTTTDMEREFMLLEQAGALLEEQVDGLLLET FYDEFELLHA
  VKVLRKQTNIPIVAQLALHEAGTTQNGNDVNEILKQFIDYGANVVGLNCQLGPLHMT EAFKMISIPQNGYLSAYPN
  AGLPNYVEGRYVYEGSPAYFEAMTPNFIEQGIRLLGGCCGTTPEHIQSMKRAVANITPVIEKETIQRPKV VHTHEK
  RSKAHVT LAEKAKKQTTVVVELDPPKTLDTQRFFEGARALKRAGADAITLADNSLASPRVSNMAMGALLTKHDIPV
  LTHLTRDHNVIQLQSHLLGLSALGMEEVLALTGD PARVGDFPGATSVYDLSSI ELIKMIKEMNDGRSILGKSIGP
  ATRFSVGGAFNPHVRHLKAAVKRMERKIDAGAEYFLTQPIYDIK LIEEVYEATKHLEQPIFIGIMPLVSKRNADFL
  HFEVPGITLPEEIRERMDGHETKEAAIEEGIRISQELVDAAMKYFNGIY LITPFLKYEITEHLVKYVREKQEVKEG
  IN
```

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