

# 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：

```
$ fastx-utils dedup

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`

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

```
>AGE80385.1
MKKIIIIISATTIVIGITSFAYFGSKTPLHNEAKAVESQKHNNHKKEEIPAFPKADHNAKKIDNDFSVVTNPKSNLVLINK
HRKLPDGYIPEDLTRPNVPFISPKDKEKTLRLKDAEALENMFKAACKKEGLDLTAVSGYRSYKRQKSLHDTYVRRQGKAE
ANSVSAIPGTSEHQTGLAMDISSSAKFQLEPIFGETAEGKWVAEHAHEFGFVIRYLEDKTDTTEYAYEPWHLRYVGPNPY
ATYLYKHHLTLEEAMEDKK
>AGE79558.1
MINHELVERINFLAKKAKAEGLTEEEQRERQSLREQYLKGFRQNMLNELKGIKVVNEEGTDVPTPKLKALKKQDNAKLN
>AGE81073.1
MDKVISNEILQQFKDRMRLGDDEDANLRRILFASNKDLTRVCGNYDLNIDEVFKELVFERSRYVYNDALEYFDKNFLSQI
NLSIGKALEAIKLDGD
>AGE80385.1
MQNGKVKWFNSEKGFGFIEVEGGEDVFVHFSAIQGEGFKTLEEGQEVTFEVEQGNRGPQATNVNKK
>AGE81522.1
MKKTILTTTAALTMMGTGMGINVDHIKPAEVKADTISFYDVPNNHWATKAITNLANRNIVVGYGNGQFGFGDNVTRGQVA
RMIYNYLKPADAGNFKNPFSDIKGHMFEKEILALAKVGIIKGYGEGKFGPDDILTREQMAQVLTNAFKFEGTKKTSFVDV
DKNSWSYKAIGALEEKGVITGTGGNMYSPTSVVTREQYSQFLFNSINVIEKETKPEEKPNTGGEVKPEEKPNTGGEVKPE
EKPNTGEETKPVNIPEWLETSLATNDFTFTQAWYDGSEAINKAASTNAQQIVKNINSKYGTNLKYSEVGAIVQLVDGARE
QLWLAGMNVNDFRVTFRVSNNAMIELTKELVTLVNSDLNLDQEIQEIIPSAPMKIKNVEKGDYKIRISPAMADQMITIIIE
KK
```

运行命令： 去除 `protein.faa` 文件中的重复序列。

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

```
>AGE80385.1
MKKIIIIISATTIVIGITSFAYFGSKTPLHNEAKAVESQKHNNHKKEEIPAFPKADHNAKKIDNDFS SVVTNPKSNLVLINK
HRKLDPDGYIPEDLTRPNVPFISPKDKEKTL LRKDAAEAL ENMFKA AKKEGLDLTAVSGYRSYKRQKSLHDTYVRRQGKAE
ANSVSAIPGTSEHQTGLAMDISSSAKFQLEPIFGETAEGKWVAEHAHEFGFVIRYLEDKTDTTEYAYEPWHLRYVGNPY
ATYLYKHHLTLEEAMEDKK
>AGE79558.1
MINHELVERINFLAKKAKAEGLTEEEQRERQSLREQYLKGFRQNMLNELKGIKVVNEEGTDVTPTKLKALKKQDNAKLN
>AGE81073.1
MDKVISNEILQQFKDRMRLGDDEDANLRRILFASNKDLTRVCGNYDLNIDEVFKELVFERSRYVYNDALEYFDKNFLSQI
NSLSIGKALEAIKLDGD
>AGE81522.1
MKKTILTTTAALTMMGTGMGINVDHIKPAEVKADTISFYDVPNNHWATKAITNLANRNIVVGYGNGQFGFGDNVTRGQVA
RMIYNYLKPADAGNFKNPFSDIKGHMFEKEILALAKVGIIKGYGEGKFGPDDILTREQMAQVLTNAFKFEGTKKTSFVDV
DKNSWSYKAIGALEEKGVITGTGGNMYSPTSVVTREQYSQFLFNSINVIEKETKPEEKPNTGGEVKPEEKPNTGGEVKPE
EKPNTGEETKPVNIPEWLETSLATNDFTFTQAWYDGSEAINKAASTNAQQIVKNINSKYGTNLKYSEVGAIQVLVDGARE
QLWLAGMNVNDFRVTFRVSNNAMIELTKELVTLVNSDLNLDQEIQEIPSAPMKIKNVEKGDYKIRISPAMADQMITIIIE
KK
>AGE80137.1
MKLLDLLSKGIVIGDGAVGTL LSHSHGLQSSFEELNVSDPDLIISIHKQYVAAGADV IQTNTYGANEA KLRMYGLENQVTK
INRAAVKLAKASVTDKNAILGTIGGMKHIGAVTTTDMEREFMLLEQAGALLEEQVDG L LLET FYDEFELLHAVKVL RKQT
NIPIVAQLALHEAGTTQNGNDVNEILKQFIDYGANVVGLNCQLGPLHMT EAFKMISIPQNGYL SAYPNAGLP NYVEGRYV
YEGSPAYFEAMTPNFIEQGIRLLGGCCGTTPEHIQSMKRAVANITPVIEKETIQRPKV VHTHEKRSKAHVTLAEKAKKQT
TVVVELDPPKTLDTQRFFEGARALKRAGADAITLADNSLASPRVSNMAMGALLTKHDIPVLTHLTCRDHNVI GLQSHLLG
LSALGMEEV LALTGD PARVGD FPGATS VYDLSSIELIKMIKEMNDGRSILGKSIGPATRFSVGGAFNPHVRHLKAAVKRM
ERKIDAGAEYFLTQPIYDIKLIEEVYEATKHLEQPIFIGIMPLVSKRNADFLHFEVPGITLPEEIRERMDGHETKEAAIE
EGIRISQELVDAAMKYFNGIY LITPFLKYEITEHLVKYVREKQEVKEGIN
```

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