

tsv-utils之多列合并: collapse

一、tsv-utils collapse介绍

功能描述:

`tsv-utils collapse` 将指定的多列元素合并为一列, 并可以指定合并元素之间的分隔符;

命令行接口:

```
1 $ tsv-utils collapse
2
3 Usage: tsv-utils collapse [options] <tsv>
4 Options:
5   -f fields to collapse. Comma separated, ie: 1,2,3.
6   -d delim.
```

可选参数:

```
1 -f 指定要合并连续的列, 如1,2,3;
2 -d 指定元素之间的分隔符;
```

二、使用场景实例及其用法(更换使用用例)

示例演示

示例文件: `megablast.txt`, `SILVA_123_SSURef.txt`

```
1 $ cat blastn.txt | head -n6
```

```
1 ZOTU_1 JQ088343.1.1479 100.000 404 0 0 1 404 355
   758 0.0 747
2 ZOTU_2 JF138741.1.1356 100.000 429 0 0 1 429 323
   751 0.0 793
3 ZOTU_3 LN568336.1.1321 98.765 405 4 1 1 404 298
   702 0.0 719
4 ZOTU_4 FJ672447.1.1441 99.764 424 1 0 1 424 355
   778 0.0 778
5 ZOTU_5 EU542474.1.1489 100.000 424 0 0 1 424 353
   776 0.0 784
6 ZOTU_6 JQ072613.1.1343 100.000 404 0 0 1 404 270
   673 0.0 747
```

```
1 $ cat SILVA_123_SSURef.txt | head -n6
```

```

1 | X92134.1.1483
   | Bacteria;Proteobacteria;Gammaproteobacteria;Pseudomonadales;Pseudomonadaceae;
   | Pseudomonas;bacterium 52N3
2 | EF442993.1.1422
   | Bacteria;Proteobacteria;Deltaproteobacteria;Desulfobacterales;Desulfobulbacea
   | e;Desulfobulbus;Desulfobulbus sp. DSM 2033
3 | JQ972882.1.1517
   | Bacteria;Actinobacteria;Actinobacteria;Streptosporangiales;Thermomonosporacea
   | e;Actinomadura;Actinomadura bangladeshensis
4 | FJ799133.1.1480
   | Bacteria;Tenericutes;Mollicutes;Acholeplasmatales;Acholeplasmataceae;Acholepl
   | asma;bacterium enrichment culture clone BA75
5 | AY554420.1.1452
   | Bacteria;Bacteroidetes;Bacteroidia;Bacteroidales;Porphyromonadaceae;Proteinip
   | hilum;Bacteroides sp. 22C
6 | AY741401.1.1512 Bacteria;Proteobacteria;Gammaproteobacteria;NKB5;Legionella-
   | like amoebal pathogen HT99

```

运行命令:

准备 `annotation.txt` 文件

```

1 | $ cut -f 1,2 megablast.txt | tsv-utils annotation -c 2 SILVA_123_SSURef.txt
   | - | tee annotation.txt | head -n6

```

```

1 | ZOTU_1 JQ088343.1.1479
   | Bacteria;Proteobacteria;Epsilonproteobacteria;Campylobacterales;Campylobacter
   | aceae;Arcobacter;uncultured bacterium
2 | ZOTU_2 JF138741.1.1356
   | Bacteria;Proteobacteria;Gammaproteobacteria;Oceanospirillales;Halomonadaceae;
   | Halomonas;uncultured bacterium
3 | ZOTU_3 LN568336.1.1321
   | Bacteria;Proteobacteria;Alphaproteobacteria;Rhodospirillales;Acetobacteraceae
   | ;Roseomonas;uncultured bacterium
4 | ZOTU_4 FJ672447.1.1441
   | Bacteria;Bacteroidetes;Bacteroidia;Bacteroidales;Porphyromonadaceae;Proteinip
   | hilum;uncultured bacterium
5 | ZOTU_5 EU542474.1.1489 Bacteria;Bacteroidetes;Bacteroidetes
   | vadinHA17;uncultured bacterium
6 | ZOTU_6 JQ072613.1.1343
   | Bacteria;Proteobacteria;Alphaproteobacteria;Rhizobiales;Rhizobiaceae;Rhizobiu
   | m;uncultured bacterium

```

合并第二列和第三列, 使用 `-f2` 参数.

```

1 | $ tsv-utils collapse -f2 -d':' annotation.txt | head -n 6

```

```
1 ZOTU_1
  JQ088343.1.1479:Bacteria;Proteobacteria;Epsilonproteobacteria;Campylobacterales;Campylobacteraceae;Arcobacter;uncultured bacterium
2 ZOTU_2
  JF138741.1.1356:Bacteria;Proteobacteria;Gammaproteobacteria;Oceanospirillales;Halomonadaceae;Halomonas;uncultured bacterium
3 ZOTU_3
  LN568336.1.1321:Bacteria;Proteobacteria;Alphaproteobacteria;Rhodospirillales;Acetobacteraceae;Roseomonas;uncultured bacterium
4 ZOTU_4
  FJ672447.1.1441:Bacteria;Bacteroidetes;Bacteroidia;Bacteroidales;Porphyromonadaceae;Proteiniphilum;uncultured bacterium
5 ZOTU_5 EU542474.1.1489:Bacteria;Bacteroidetes;Bacteroidetes vadinHA17;uncultured bacterium
6 ZOTU_6
  JQ072613.1.1343:Bacteria;Proteobacteria;Alphaproteobacteria;Rhizobiales;Rhizobiaceae;Rhizobium;uncultured bacterium
```

注意事项: 当前实现只支持单字符分隔符.

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: 8/30/2020 7:21:56 PM