

实验 03 过滤器与正则表达式

班级：数据科学与大数据 2 班

学号：202026203039

姓名：赖丽婷

用户名：llt

一、实验目的

1. 练习操作文件和目录的常用命令

二、实验要求

1. 填写实验报告，请将关键命令及其结果进行截图(确保截图中的文字清晰可见)
2. 导出为 pdf 文件，文件名为用户名-姓名-lab02.pdf，在规定截止时间之前上传作业)
3. 以下步骤中所有 s01 请换成你自己的用户名。

三、实验步骤

- 1.对实验文件 sea.txt 中所有出现的单词进行计数，并按照出现次数从大到小的顺序打印出各单词及其出现次数。(要求统一大小写并且把标点符号和多余的空白符去除干净后进行统计)

```
1 youve
llt@llt-virtual-machine:~/下载$ cat sea.txt | tr 'A-Z' 'a-z' | tr -s ' ' | tr '
' '\n' | tr -d '[:punct:]' | sort | uniq -c | sort -k1nr
2315 the
1259 and
1143 he
541 of
490 it
465 i
455 to
446 his
435 was
397 a
359 in
295 that
271 fish
248 old
215
```

cat 读取文件 把大写转为小写 删除连续的空格 把空格转为换行符 删除特殊字符 先排序 为了能删除重复的行 然后在按字符出现的数字排序

2. 打印实验文件 rfile.txt 中包含合法数字（整数、小数、科学计数法，可正可负）的行示例：

123 # 合法数字

12.3 # 合法数字

.123 # 合法数字

-3.14 # 合法数字

-3.14e9 # 合法数字

-.14E10 # 合法数字

-.14E-10 # 合法数字

1.2.3 # 非法数字

1.2.e3.5 # 非法数字

-1.-2E5 # 非法数字

4. 打印实验文件 `rfile.txt` 中包含合法身份证号（15 位或 18 位）的行
身份证编码规则：

规则 1：15 位身份证号

- 6 位地址码：任意数字
- 6 位出生日期：2 位年+2 位月(01-12)+2 位日(01-31){}
- 3 位顺序码：任意数字

规则 2：18 位身份证号

- 6 位地址码：任意数字
- 8 位出生日期：
- 4 位年(18XX,19XX,2XXX)+2 位月(01-12)+2 位日(01-31)
- 3 位顺序码：任意数字
- 1 位校验码：任意数字或大写字母 X

示例：

283987983893237 # 非法

398478348957389459 # 非法

49583749573593593845 # 非法

34985739485793578X # 非法

34895738953498934x # 非法

209349891231838 # 合法

897345151301381 # 非法

084594991042948 # 非法

094589200109193982 # 合法

95680920000139568X # 非法

39485719991016487X # 合法

```
39485719991016487X
llt@llt-virtual-machine:~/下载$ grep -E "^[0-9]{6}[0-9]{2}([0-9]{1}|1[0-2]{1})
([0-9]{1}|3[0-1]{1}|[1-2]{1}[0-9]{1})[0-9]{3}$|^[0-9]{6}([18|19|2[0-9]{1})[0-9]
{2}([0-9]{1}|1[0-2]{1})([0-9]{1}|3[0-1]{1}|[1-2]{1}[0-9]{1})[0-9]{3}[0-9|X]$"
rfile.txt
209349891231838
094589200109193982
39485719991016487X
llt@llt-virtual-machine:~/下载$
```