

练习参考

IP地址

匹配合法的IP地址

```
1 192.168.1.150
2 0.0.0.0
3 255.255.255.255
4 17.16.52.100
5 172.16.0.100
6 400.400.999.888
7 001.022.003.000
8 257.257.255.256
```

```
(\d{1,3}\.){3}\d{1,3}
```

```
(?:\d{1,3}\.){3}\d{1,3} # 400.400.999.888
```

对于ip地址验证的问题

- 可以把数据提出来后，交给IP地址解析库 处理，如果解析异常，就说明有问题，正则的验证只是一个初步的筛选，把明显错误过滤掉。
- 可以使用复杂的正则表达式验证地址正确性
- 前导0是可以的

```
1 import socket
2 nw = socket.inet_aton('192.168.05.001') # 错了抛异常
3 print(nw, socket.inet_ntoa(nw))
```

```
1 分析：
2 每一段上可以写的数字有1、01、001、000、23、023、230、100，也就说1位就是0-9，2位每一位也是0-9，3位第一位只能0-2，其余2位都可以0-9
3 (?:([0-2]\d{2}|\d{1,2})\.){3}([0-2]\d{2}|\d{1,2}) # 解决超出200的问题，但是256呢？
4 200是特殊的，要再单独分情况处理
5 25[0-5]|2[0-4]\d|[01]?\d\d? 这就是每一段的逻辑
6 (?: (25[0-5]|2[0-4]\d|[01]?\d\d?)\. ){3}(25[0-5]|2[0-4]\d|[01]?\d\d?)
```

提取文件名

选出含有ftp的链接，且文件类型是gz或者xz的文件名

```

1 ftp://ftp.astron.com/pub/file/file-5.14.tar.gz
2 ftp://ftp.gmp1ib.org/pub/gmp-5.1.2/gmp-5.1.2.tar.xz
3 ftp://ftp.vim.org/pub/vim/unix/vim-7.3.tar.bz2
4 http://anduin.linuxfromscratch.org/sources/LFS/lfs-
  packages/conglomeration//iana-etc/iana-etc-2.30.tar.bz2
5 http://anduin.linuxfromscratch.org/sources/other/udev-lfs-205-1.tar.bz2
6 http://download.savannah.gnu.org/releases/libpipeline/libpipeline-
  1.2.4.tar.gz
7 http://download.savannah.gnu.org/releases/man-db/man-db-2.6.5.tar.xz
8 http://download.savannah.gnu.org/releases/sysvinit/sysvinit-2.88dsf.tar.bz2
9 http://ftp.altlinux.org/pub/people/legion/kbd/kbd-1.15.5.tar.gz
10 http://mirror.hust.edu.cn/gnu/autoconf/autoconf-2.69.tar.xz
11 http://mirror.hust.edu.cn/gnu/automake/automake-1.14.tar.xz

```

```

1 .*ftp.*\.(?:gz|xz)
2 .*ftp.*(/[^\.]+(?:.xz|.gz))$ # 后面的分组就是文件名
3 (?:<=.*ftp.*\/)[^\.]*\.(?:gz|xz) # 断言文件名前一定还有ftp, python的re模块不支持断言中
  使用*、+或{3,4}等不固定长度的匹配

```

匹配邮箱地址

```

1 test@hot-mail.com
2 v-ip@magedu.com
3 web.manager@magedu.com.cn
4 super.user@google.com
5 a@w-a-com

```

确定规则, 邮箱名为字母、数字、下划线、减号、点号, 域名一样

```
1 \w[\w.-]*@([\w-]+\.)+\w+
```

@前至少一个字符开头, 可以有字母、数字、下划线、减号、点号

@后至少有一个字符后跟点, 重复至少1次, 最后以1个字符串结尾

匹配html标记

提取href中的链接url, 提取文字“马哥教育”

```
1 <a href='http://www.magedu.com/index.html' target='_blank'>马哥教育</a>
```

有的时候可以这样写 `马哥教育`, 属性attr不标准的写法不写单双引号, 这个一般交给浏览器可以正常处理。

```
1 href=['"]?([^"'\s]*)
```

提取内容时, 要考虑一下情况

```
1 <a href='http://www.magedu.com/index.html' target='_blank'>马哥教育</a><a>马哥
  教育1</a>
```

```
1 <a.*?>[<>]*
```

a标签要到第一个右尖括号结束，不能贪婪。内容是其之后第一个左尖括号之前

匹配URL

```
1 http://www.magedu.com/index.html
2 https://login.magedu.com
3 file:///ect/sysconfig/network
```

```
1 (\w+):\/\/(\S*)
```

匹配二代中国身份证ID

```
1 321105700101003
2 321105197001010030
3 11210020170101054X
4 17位数字+1位校验码组成
5 前6位地址码，8位出生年月，3位数字，1位校验位（0-9或X）
```

```
1 身份证验证
2 身份证验证需要使用公式计算，最严格的应该实名验证。
3 \d{17}[0-9XX]|\d{15}
```

