# The Double Split Pattern

Sometimes we split a line one way, and then grab one of the pieces of the line and split that piece again

```
From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008
```

```
words = line.split()
email = words[1]
pieces = email.split('@')
print(pieces[1])
```

stephen.marquard@uct.ac.za

['stephen.marquard', 'uct.ac.za']

'uct.ac.za'

# The Regex Version

```
From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008

import re
lin = 'From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008'
y = re.findall('@([^ ]*)',lin)
print(y)


['uct.ac.za']
```

`'@([^ ]*)'`

Look through the string until you find an at sign

# The Regex Version

```
From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008

import re
lin = 'From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008'
y = re.findall('@([^ ]*)',lin)
print(y)

['uct.ac.za']
```

'@([^ ]*)'

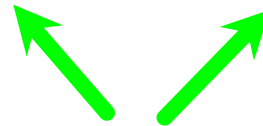Match non-blank character    Match many of them

# The Regex Version

```
From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008

import re
lin = 'From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008'
y = re.findall('@([^ ]*)',lin)
print(y)

['uct.ac.za']
```

`'@([^ ]*)'`
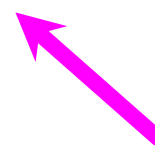
Extract the non-blank characters

# Even Cooler Regex Version

```
From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008

import re
lin = 'From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008'
y = re.findall('^From .*@([^ ]*)',lin)
print(y)

['uct.ac.za']
```

`'^From .*@([^ ]*)'`

Starting at the beginning of the line, look for the string 'From '

# Even Cooler Regex Version

```
From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008

import re
lin = 'From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008'
y = re.findall('^From .*@([^ ]*)',lin)
print(y)

['uct.ac.za']
```

### '^From .*@([^ ]*)'

Skip a bunch of characters, looking for an at sign

# Even Cooler Regex Version

From stephen.marquard@uct.ac.za Sat Jan   5 09:14:16 2008

```
import re
lin = 'From stephen.marquard@uct.ac.za Sat Jan   5 09:14:16 2008'
y = re.findall('^From .*@([^ ]*)',lin)
print(y)
```

['uct.ac.za']

**'^From .*@([^ ]*)'**

Start extracting

# Even Cooler Regex Version

**From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008**

```
import re
lin = 'From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008'
y = re.findall('^From .*@([^ ]*)',lin)
print(y)
```

['uct.ac.za']

## '^From .*@([^ ]+)'

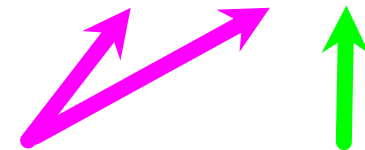Match non-blank character   Match many of them

# Even Cooler Regex Version

```
From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008

import re
lin = 'From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008'
y = re.findall('^From .*@([^ ]*)',lin)
print(y)

['uct.ac.za']
```

`'^From .*@([^ ]+)'`

Stop extracting

# Spam Confidence

```python
import re
hand = open('mbox-short.txt')
numlist = list()
for line in hand:
    line = line.rstrip()
    stuff = re.findall('^X-DSPAM-Confidence: ([0-9.]+)', line)
    if len(stuff) != 1 :  continue
    num = float(stuff[0])
    numlist.append(num)
print('Maximum:', max(numlist))
```

```
python ds.py
Maximum: 0.9907
```

```
X-DSPAM-Confidence: 0.8475
```

# Escape Character

If you want a special regular expression character to just behave normally (most of the time) you prefix it with '\'

```
>>> import re
>>> x = 'We just received $10.00 for cookies.'
>>> y = re.findall('\$[0-9.]+',x)
>>> print(y)
['$10.00']
```

At least one
or more

\$[0-9.]+

A real dollar sign        A digit or period