

# Регулярные выражения

Как уже говорилось при поиске, копировании, перемещении, etc можно использовать маски. Например гифки для своих статей я делаю через утилиту "convert". Если мне нужно сконвертировать файлы из директории "frames", в которой лежат файлы: pic1.png, pic2.png, pic3.png, pic4.png, pic5.png, pic.png, 1.png, file.txt, и мне нужны только picN.png где N от 1 до 4, то:

```
$ convert -delay 200 -loop 0 frames/pic[1..4].png test.gif
```

Вот только этого не всегда достаточно, если бы один из файлов назывался бы pic4.png, то он не был бы включен в результирующую гифку или не был бы скопирован или найден и etc. Когда требуется гибкость, то используются регулярные выражения.

## Регулярные выражения

Регулярные выражения - формальный язык поиска и осуществления манипуляций с подстроками в тексте, основанный на использовании метасимволов. Для поиска используется строка-образец (*pattern*), состоящая из символов и метасимволов и задающая правило поиска. Для манипуляций с текстом дополнительно задается строка замены, которая также может содержать в себе специальные символы. (Wiki)

### Синтаксис регулярных выражений

Регулярки (регекспы, regex, regular expression) являются регистрозависимыми. Строка "hello" - это уже регулярка, которой соответствует только одна строка - "hello" и не соответствуют: "Hello", "helLo", etc. Также при составлении регулярного выражения применяются спецсимволы (или метасимволы): `.` `^` `$` `*` `+` `?` `{}` `[]` `\` `|` `()` `<` `>` `:`. Поэтому если требуется найти спец символ, то он должен быть экранирован с помощью `"\"`. А еще есть предопределенные соответствия, что-то типо шоткатов:

- **\d** - Соответствует любой цифре; эквивалент класса **[0-9]**.
- **\D** - Соответствует любому нечисловому символу; эквивалент класса **[^0-9]**.
- **\s** - Соответствует любому символу whitespace; эквивалент **[\t\n\r\f\v]**.
- **\S** - Соответствует любому не-whitespace символу; эквивалент **[^\t\n\r\f\v]**.
- **\w** - Соответствует любой букве или цифре; эквивалент **[a-zA-Z0-9\_]**.
- **\W** - Наоборот; эквивалент **[^a-zA-Z0-9\_]**.

Итак, начнем с простого, с перечислений. Перечисления заключаются в **[]** между которыми указывается набор символов один из которых должен быть в строке. Например в результате команды "split" мы получаем файлы с меняющимися индексами: chunkaa, chunkab, etc. Таким образом (если у нас не больше 28-ми кусков) мы можем найти их все с помощью такого регулярного выражения: chunka**[abcdefghi.....короче весь алфавит]**. Если у тебя хватит терпения закончить весь алфавит, то получится то что нужно. Таким образом процессор регулярных выражений будет искать все, что начинается со строки "chunka" + один символ из набора. Правда когда нужно указать алфавит, то можно поступить иначе, указать диапазон: **a-z**. Получается: chunka**[a-z]**. Как просто, правда? А если в последней букве может быть другой регистр? Тогда: chunka**[A-z]**, что будет равно chunka**[a-zA-Z]**. Тоже самое и с цифрами: **[0-9]**. Еще можно указывать не полный диапазон, а его срез, например **[1-4]**. Также в наборе может быть несколько диапазонов и перечислений:

chunka[\[a-cDg0-3\]](#). Т.е. на месте последней буквы нашего регулярного выражения может быть любой символ из [\[abcDg0123\]](#) (порядок, кстати, не важен).

Открой сайт [RegExr](#) и пробуй там. Ты можешь заметить в строке ввода регулярки на сайте в начале и в конце `"/"`. Это потому что регулярные выражения пришли в наш мир из языка программирования Perl, а в перле начало и конец регулярного выражения обозначается `"/"`, а механизм поиска с помощью регулярок встроен в язык. Есть еще рожденный в грехе язык программирования PHP в котором при написании регекспа нужно на концах указать любые спец символы которые не будут использованы в самом выражении идиотизм.

Описание значений метасимволов:

- `.` - обозначает любой символ (кроме символа перевода строки!).
- `^` - указатель на начало потока данных, в `[]` используется для исключения символов из поиска.
- `$` - указатель на конец потока данных.
- `*` - повторяющий символ, означающий ничего или много.
- `+` - 1 или много.
- `?` - 0 или 1 и определяет жадность.
- `{ }` - определяет диапазон повторений.
- `[]` - указывает набор.
- `\` - экранирование.
- `|` - логическое ИЛИ.
- `()` - неименованная группировка.
- `<>` - именнованная группировка.

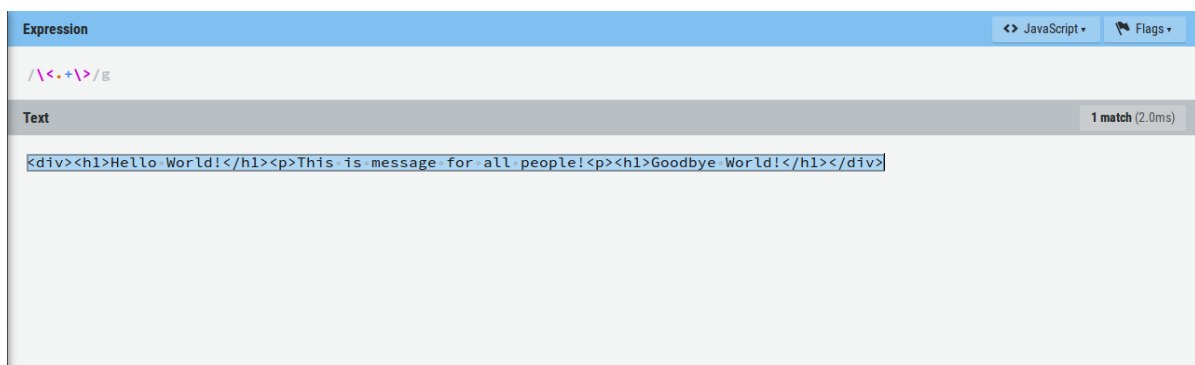
Регулярные выражения могут быть жадными и нежадными.

```
<div><h1>Hello World!</h1><p>This is message for all people!<p></div>
```

Попробуем найти все теги в этом тексте. Если мы напишем такую регулярку

```
<.+>
```

то она съест все от первого символа `"<"` до последнего `">"`, это и есть жадное регулярное выражение. Символ `">"` ищется самый последний.



Чтобы умерить аппетиты регулярки используется символ `"?"`. Т.е. чтобы найти именно теги надо после `"+"` написать `"?"`

```
<.+?>
```

```
Expression: /\<.+?\/\>/g
Text: <div><h1>Hello World!</h1><p>This is message for all people!<p><h1>Goodbye World!</h1></div>
8 matches (2.0ms)
```

Отлично! Вот только не все так просто, если кто то решил написать в тексте "<>", то получится следующее

```
Expression: /\<.+?\/\>/g
Text: <div><h1>Hello <> World!</h1><p>This is message for all people!<p><h1>Goodbye World!</h1></div>
8 matches (0.0ms)
```

Потому что ".+" говорит, что любой символ, а символ ">" тоже к ним относится, должен повторяться 1 или более раз. Таким образом строка "<> Any text>" попадет под наше условие из-за того, что между первыми скобками нет ничего. Чтобы решить это нужно заменить символ повторения на "\*".

```
Expression: /\<.*?\/\>/g
Text: <div><h1>Hello <> World!</h1><p>This is message for all people!<p><h1>Goodbye World!</h1></div>
9 matches (2.0ms)
```

Хорошо, а что если тебе понадобится найти текст между тегами? Тогда регулярное выражение может быть таким:

`<[w\s!]+<`

Тут жадность не важна потому, что мы задали набор из букв, пробела и восклицательного знака.

Хорошо. А что если мы ищем строки только определенной длины? Например длина строки может быть от 4-х до 5-ти символов. Тогда нужно указать интервал в {} через запятую.

`\w{4,5}`

А если нужны слова?

`[^\w]\w{4,5}\s`

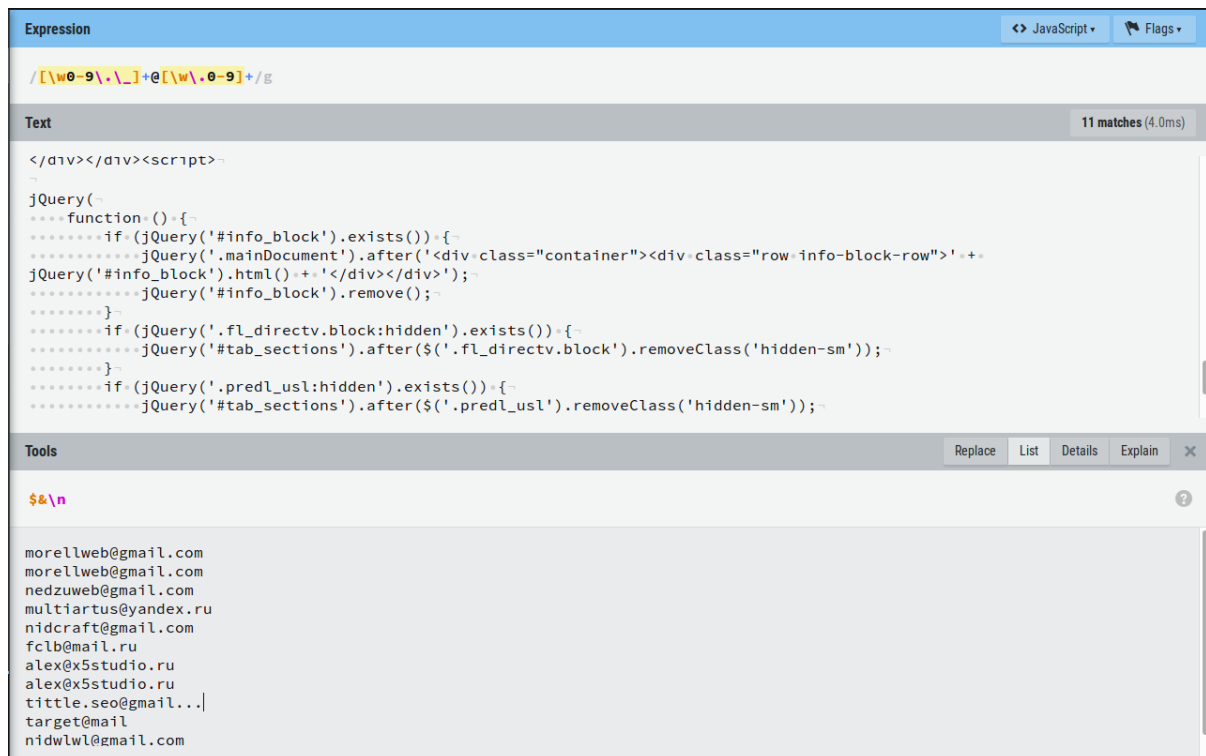
Также кол-во повторений можно задавать так: `{4}` - равно 4-м, `{4,}` - от 4-х и больше, `{,4}` - не больше 4-х.

Пример посерьезнее, возьмем главную страницу сайта "<https://freelance.ru/>" и попробуем найти там email адреса. Исходный код можно посмотреть через ПКМ -> Исходный код страницы.

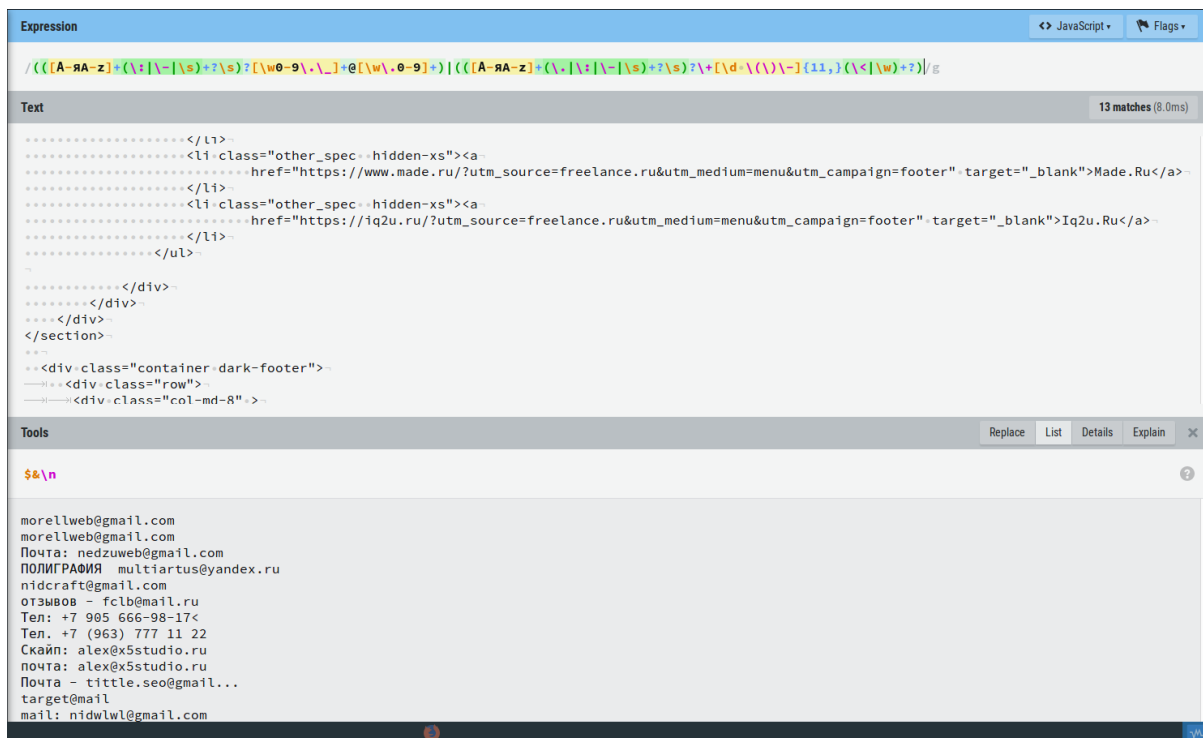
Пишем регулярку! Мы знаем, что в email адресе в качестве имени пользователя могут быть: цифры, буквы, знак подчеркивания и точка; все это может повторяться. В качестве разделителя идет символ "@". В домене могут быть: цифры, буквы, точка. Тогда регулярка будет такой:

`[w0-9\.\_ ]+@[w\.\_0-9]+`

И результат:



Также можно собрать мета информацию о контакте и искать номера телефонов, например.



Разберемся что я тут написал:

`((([A-ЯА-Я-z]+(\.|\-|\s)?\s)?[w0-9\.\ ]+@[w\.\0-9]+)|((([A-ЯА-Я-z]+(\.|\-|\s)?\s)?\+[d\(\)\-]{11,})(\<|w)+?))`

Делится все по скобкам и получается что у нас 2 регулярных выражения объединенных с помощью "|" в один:

`([A-ЯА-Я-z]+(\.|\-|\s)?\s)?[w0-9\.\ ]+@[w\.\0-9]+`

`([A-ЯА-Я-z]+(\.|\-|\s)?\s)?\+[d\(\)\-]{11,}(\<|w)+?`

В первом и втором регулярках тоже есть группировка, первая часть обеих регулярок повторяется:

`([A-ЯА-Я-z]+(\.|\-|\s)?\s)?`

Эта часть заключена в скобки чтобы в конце поставить знак "?" (это называется *группированием*). И работать это будет так, если соответствие этому регулярному выражению найдено, то оно берется в результат, иначе просто игнорируется.

Расчленим его и посмотрим на куски:

- **[A-ЯА-Я-z]+** - Будет сопоставляться кириллица и латинница, "+" - так как должна быть либо буква либо несколько букв которые идут до выражения ниже
- **(\.|\-|\s)+?\s** - здесь объединяется группа символов разделителей, ведь описание контакта как то отделяется от самого контакта. Это может быть один символ из: `.\- \s`. Вконец стоит "+", что значит символы могут повториться и "?" чтобы не жадничал. завершается все символом "s". Все это нужно чтобы соответствовать: Тел. {data}, Email - {data}, Email: {data}.

С первой частью разобрались. Теперь разберемся со второй частью второй регулярки:

`\+[d\(\)\-]{11,}(\<|w)+?`

Глаза можно выколоть смотря и составляя такие сложные регулярные выражения, да? Тогда вот посмотри на регулярное выражение проверяющее валидность email адреса по стандартам RFC.

[illegible]

[illegible]

]])\*)\*@(?(?:\r\n)?[ \t])\*(?:^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t  
])+\|Z|(?=[\"]())<>@,;:\\".\\[ ]))\|([^\[\]\r\\]\\".)\*\](?:?:\r\n)?[ \t])\*(?  
:\\".?(?:\r\n)?[ \t])\*(?:^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t])+\|  
Z|(?=[\"]())<>@,;:\\".\\[ ]))\|([^\[\]\r\\]\\".)\*\](?:?:\r\n)?[ \t])\*)\*(?:  
^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t])+\|Z|(?=[\"]())<>@,;:\\".\\  
[ ]))\"(?:^[^\r\\]\\".|(?:?:\r\n)?[ \t]))\"(?:?:\r\n)?[ \t])\*)\*<(?:?:\r\n)  
?[ \t])\*(?:@:(?:^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t])+\|Z|(?=[\"]  
())<>@,;:\\".\\[ ]))\|([^\[\]\r\\]\\".)\*\](?:?:\r\n)?[ \t])\*(?:\\".?(?:\r\n)  
?[ \t])\*(?:^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t])+\|Z|(?=[\"]())<>  
@,;:\\".\\[ ]))\|([^\[\]\r\\]\\".)\*\](?:?:\r\n)?[ \t])\*)\*(?:,@(?:?:\r\n)?[  
 \t])\*(?:^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t])+\|Z|(?=[\"]())<>@,  
;:\\".\\[ ]))\|([^\[\]\r\\]\\".)\*\](?:?:\r\n)?[ \t])\*(?:\\".?(?:\r\n)?[ \t]  
)\*(?:^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t])+\|Z|(?=[\"]())<>@,;:\\"  
.\\[ ]))\|([^\[\]\r\\]\\".)\*\](?:?:\r\n)?[ \t])\*)\*)\*(?:?:\r\n)?[ \t])\*)?  
(?:^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t])+\|Z|(?=[\"]())<>@,;:\\".  
\\[ ]))\"(?:^[^\r\\]\\".|(?:?:\r\n)?[ \t]))\"(?:?:\r\n)?[ \t])\*(?:\\".?(?::  
 \r\n)?[ \t])\*(?:^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t])+\|Z|(?=[\"]  
())<>@,;:\\".\\[ ]))\"(?:^[^\r\\]\\".|(?:?:\r\n)?[ \t]))\"(?:?:\r\n)?[ \t])  
)\*)\*@(?:?:\r\n)?[ \t])\*(?:^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t])  
+\|Z|(?=[\"]())<>@,;:\\".\\[ ]))\|([^\[\]\r\\]\\".)\*\](?:?:\r\n)?[ \t])\*(?:\  
:\\".?(?:\r\n)?[ \t])\*(?:^(<>@,;:\\".\\[ \000-\031]+(?:?:(?:\r\n)?[ \t])+\|Z  
|(?=[\"]())<>@,;:\\".\\[ ]))\|([^\[\]\r\\]\\".)\*\](?:?:\r\n)?[ \t])\*)\*>(?::  
?:\r\n)?[ \t])\*)\*)?;\s\*)

## **RFC 5322**

\A(?:[a-z0-9!#\$%&'\*/=?^\_`{}~-]+(?:\[a-z0-9!#\$%&'\*/=?^\_`{}~-]+)\*  
| "(?:\[x01-\x08\x0b\x0c\x0e-\x1f\x21\x23-\x5b\x5d-\x7f]  
| [\\x01-\x09\x0b\x0c\x0e-\x7f](#)))\*  
@ (?:?:[a-z0-9](?:[a-z0-9]\*[a-z0-9])?\\".)+[a-z0-9](?:[a-z0-9]\*[a-z0-9])?  
| \[(?:?:25[0-5][2[0-4][0-9][01]?[0-9][0-9]?)\\".){3}  
(?:25[0-5][2[0-4][0-9][01]?[0-9][0-9]?|[a-z0-9]\*[a-z0-9]:  
(?:\[x01-\x08\x0b\x0c\x0e-\x1f\x21-\x5a\x53-\x7f]  
| [\\x01-\x09\x0b\x0c\x0e-\x7f](#)))+)



\j)\z

## Применение регулярных выражений

Теория без практики мертва. В линукс есть замечательный инструмент "egrep", который позволяет искать файлы используя синтаксис регулярных выражений. Искать можно в файлах (*также рекурсивно*) и передавая информацию на стандартный поток ввода. Крайне удобно применять egrep при поиске процессов. Команда "ps" выводит информацию о процессах в системе.

```
$ ps aux
```

Процессов очень много, тогда если нас интересуют все процессы использующие Python можно сделать так:

```
$ ps aux | egrep python
```

Стоит учесть что egrep будет искать соответствие регулярному выражению во всей строке это значит, что процессы

```
user      5716 0.0  0.0 197212      8 ?        S    date   time /usr/bin/python3
/path/to/file

user      5786 0.0  0.0 433096 4132 ?        Sl   date   time /usr/bin/env python3

user      576  0.0  0.0 135096 132 ?        Sl   date   time /usr/bin/nano
/path/to/python3/file.txt
```

будут тоже показаны в выводе. Чтобы это исправить читайте справку по "ps" или составляйте более точные регулярки.

Есть одна особенность использования grep (*egrep аналогичная утилита только использует расширенный синтаксис по умолчанию запускаясь с флагом -e*) зависимость от системной локали. Например на моей системе те регулярки, что я писал под код сайта флиланса уже не работают. Поэтому научимся применять пайтон при решении подобных задач.

```
unit@unit-computer: /tmp/perms [12:59:59]
> $ egrep "(([A-z]+(\:|\-|\s)+?\s)?[\w0-9\.\_]+@[ \w\.\0-9]+)|(([A-z]+(\.|\:|\-|\s)+?\s)
)?\+[\d \(\)\-]{11,}(\<|\w)+?" freelance.html
grep: Неверный конец диапазона

unit@unit-computer: /tmp/perms [13:00:13]
> $
```

## Python regex

Python - это высокоуровневый интерпретируемый язык программирования общего назначения. Большого знать сейчас и не надо. По умолчанию присутствует в подавляющем числе дистрибутивов линукс, исключение составляют сборки под микроконтроллеры разве что, например его может не быть в роутерах. На данный момент дистрибутивы для десктопа распространяются с двумя версиями этого языка: 2.7 и 3.9.

```
unit@unit-computer: /tmp/perms [12:39:49]
> $ python --version
Python 2.7.12

unit@unit-computer: /tmp/perms [12:43:43]
> $ python3 --version
Python 3.5.2

unit@unit-computer: /tmp/perms [12:43:49]
> $
```

В качестве редактора я использую SublimeText 3, ты можешь пользоваться любым текстовым редактором, хоть nano.

Вот такой простой скрипт:

```
/tmp/perms/pyregex.py - Sublime Text (UNREGISTERED)
freelance.html x pyregex.py
1 #!/usr/bin/env python3
2 import sys # нужно чтобы получить параметры из терминала
3 import re # модуль для работы с регулярками
4
5 pattern = sys.argv[1] # первый параметр будет паттерном
6 target_file = sys.argv[2] # второй - файлом
7
8 result = re.findall( # запускаем поиск, для этого нужны 2 параметра
9     pattern, # подставляем паттерн
10     open(target_file, 'r').read() # читаем файл
11 )
12
13 if pattern: # если что то нашли то
14     print( # отправим результат поиска в терминал
15         *(max(x) for x in result), # что нашли (тут происходит магия)
16         sep='\n' # укажем разделитель
17     )
18 else: # иначе
19     print('Not found!') # просто скажем что ничего не найдено
20
```

И запустим:

```
unit@unit-computer: /tmp/perms [13:01:46]
> $ ./pyregex.py "([A-яA-Z]+(\:|\-|\s)+?\s)?[\w0-9\.\_]+@[ \w\.\0-9]+)|([A-яA-Z]+(\:|\-|\s)+?\s)?\+[ \d \(\)\-]{11,}(\<|\w)+?" freelance.html
morellweb@gmail.com
morellweb@gmail.com
Почта: nedzuweb@gmail.com
ПОЛИГРАФИЯ multiartus@yandex.ru
nidcraft@gmail.com
отзывов - fclb@mail.ru
Тел: +7 905 666-98-17<
Тел. +7 (963) 777 11 22 С
кайп: alex@x5studio.ru
почта: alex@x5studio.ru
Почта - tittle.seo@gmail...
target@mail
mail: nidwllwl@gmail.com

unit@unit-computer /tmp/perms [13:01:54]
> $
```

Ничего сложного :-)

Немного о синтаксисе `./file`. Думаю тебе понятно, что это указание полного пути к файлу. Первая строка в самом файле `#!/usr/bin/env python3` определяет, что будет запущено при таком подходе. Т.е. будет запущена утилита `env` которая в качестве параметра получит `"python3"` и запустит его. Это нужно чтобы не писать полный путь к интерпретатору, который может лежать где угодно. `./file` это способ относительного запуска исполняемых файлов. Думаю объяснения уже излишни.

**Порой лучший способ что то понять, взять и попробовать!**