

Welcome

INTERMEDIATE REGULAR EXPRESSIONS IN R



Angelo Zehr
Data Journalist

Where you might have left off



From Rebus to writing custom expressions

Does `"cat"` start with `"c"` ?

The rebus way:

```
str_detect("cat", pattern = START %R% "c")
```

Regular expression:

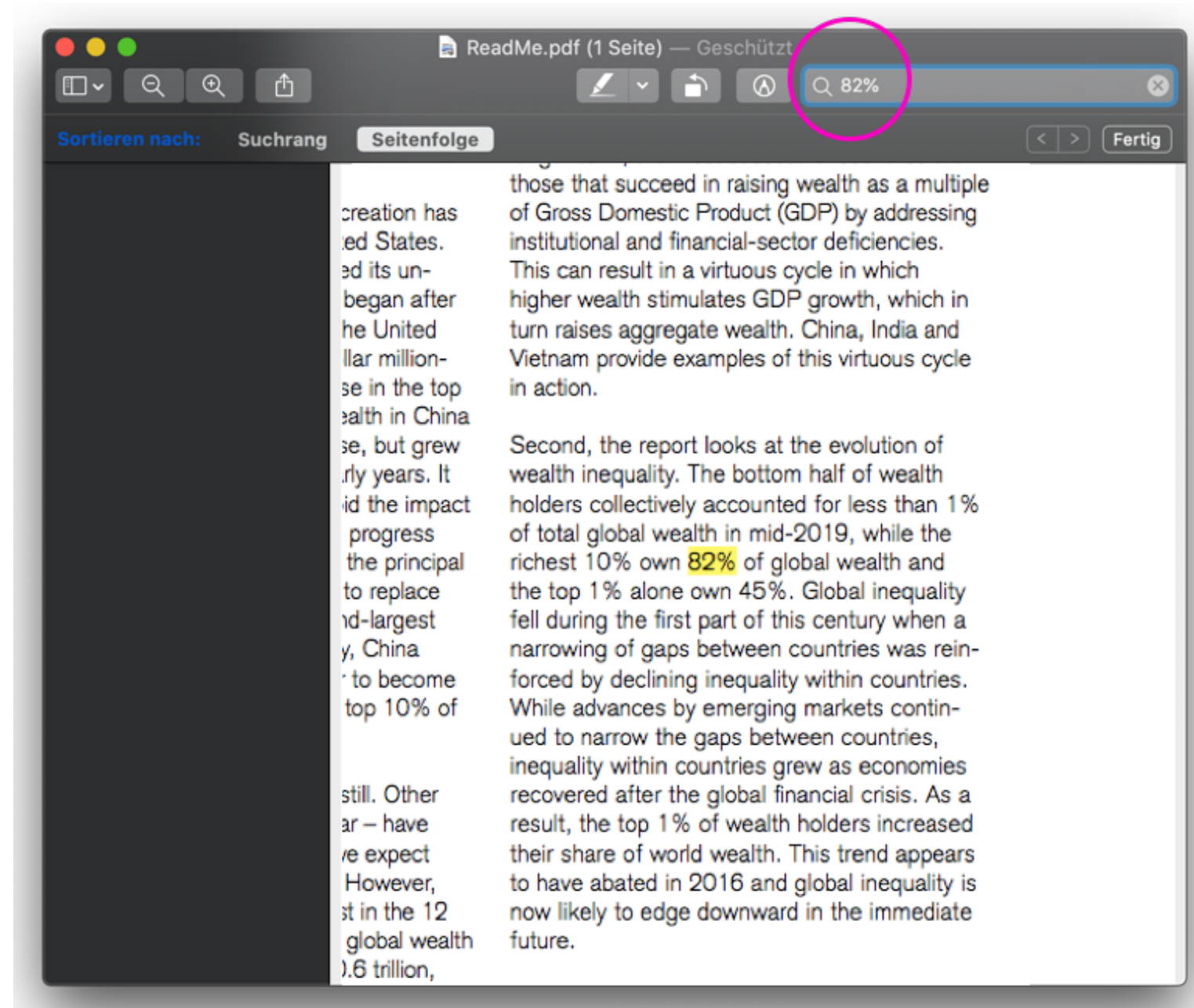
```
str_detect("cat", pattern = "^c")
```

Prerequisites: stringr

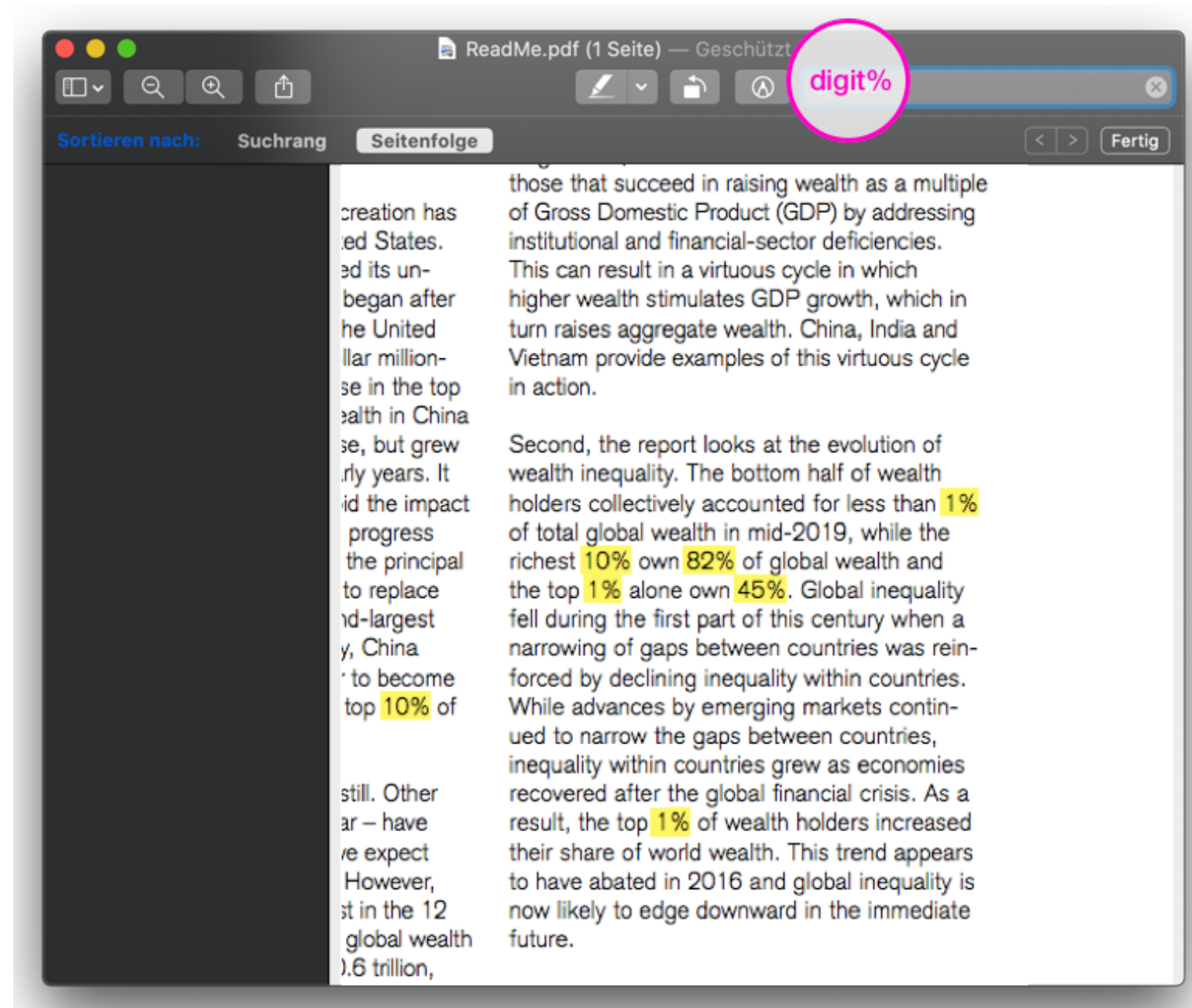
```
str_detect(string, pattern)
```

```
str_match(string, pattern)
```

What regular expressions will help you achieve



What regular expressions will help you achieve



Our first dataset

```
movie_titles <- c(  
  "Karate Kid",  
  "The Twilight Saga: Eclipse",  
  "Knight & Day",  
  "Shrek Forever After (3D)",  
  "Marmaduke.",  
  "Predators",  
  "StreetDance (3D)",  
  "Robin Hood",  
  "Micmacs A Tire-Larigot",  
  "Sex And the City 2",  
  ...  
)
```

```
movie_titles[  
  str_detect(  
    movie_titles,  
    pattern = "^K"  
  )  
]
```

```
"Karate Kid",  
"Knight & Day",  
...  
)
```

Special characters in regular expressions

Special character	Meaning
<code>^</code>	Caret: Marks the beginning of a line or string
<code>\$</code>	Dollar Sign: Marks the end of a line or string
<code>.</code>	Period: Matches anything: letters, numbers or white spaces
<code>\\.</code>	Two backslashes: Escapes the period when we search an actual period

For example

Code	Result
<code>str_match("Book", "^.")</code>	Will match <code>"B"</code>
<code>str_match("Book", ".\$")</code>	Will match <code>"k"</code>
<code>str_match("Book", "\\.")</code>	No match
<code>str_match("Book.", "\\.")</code>	Will match <code>"."</code>

Let's practice!

INTERMEDIATE REGULAR EXPRESSIONS IN R

Character classes and repetitions

INTERMEDIATE REGULAR EXPRESSIONS IN R



Angelo Zehr
Data Journalist

Available character classes

Character Class	Example
<code>\\d</code> or <code>[:digit:]</code>	0, 1, 2, 3, ...
<code>\\w</code> or <code>[:word:]</code>	a, b, c..., 1, 2, 3..., _
<code>[A-Za-z]</code> or <code>[:alpha:]</code>	A, B, C, ..., a, b, c, ...
<code>[aeiou]</code>	either a , e , i , o or u
<code>\\s</code> or <code>[:space:]</code>	" " , tabs or line breaks

A concrete example

<code>str_match_all()</code>	Result
<code>"Hi John_35", "\\d"</code>	<code>"3"</code> , <code>"5"</code>
<code>"Hi John_35", "\\w"</code>	<code>"H"</code> , <code>"i"</code> , <code>"J"</code> , <code>"o"</code> , <code>"h"</code> , <code>"n"</code> , <code>"_"</code> , <code>"3"</code> , <code>"5"</code>
<code>"Hi John_35", "[A-Za-z]"</code>	<code>"H"</code> , <code>"i"</code> , <code>"J"</code> , <code>"o"</code> , <code>"h"</code> , <code>"n"</code>
<code>"Hi John_35", "[aeiou]"</code>	<code>"i"</code> , <code>"o"</code>
<code>"Hi John_35", "\\s"</code>	<code>" "</code>

Repetitions

Syntax	Meaning
<code>\\w{2}</code>	exactly 2 times
<code>\\w{2, 3}</code>	minimum 2 times, maximum 3 times
<code>\\w{2, }</code>	minimum 2 times, but no maximum
<code>\\w+</code>	1 or more repetitions
<code>\\w*</code>	0, 1 or more repetitions

Inversion of character classes

Original	Negation
<code>\\d</code> match digits	<code>\\D</code> match all but digits
<code>\\w</code> match word characters	<code>\\W</code> match all but word characters
<code>\\s</code> match spaces	<code>\\S</code> match all but spaces
<code>[a-zA-Z]</code> match alphabet	<code>[^a-zA-Z]</code> match all but alphabet

Custom pattern with classes

```
str_match_all("Toy Story 3", "[\\d\\s]")
```

Result:

```
      [,1]  
[1,] " "  
[2,] " "  
[3,] "3"
```


Let's practice!

INTERMEDIATE REGULAR EXPRESSIONS IN R

The pipe and the question mark

INTERMEDIATE REGULAR EXPRESSIONS IN R



Angelo Zehr
Instructor

This or that

```
lines <- c(
  "Karate Kid 2, Distributor: Columbia, 58 Screens",
  "Finding Nemo, Distributors: Pixar and Disney, 10 Screens",
  "Finding Harmony, Distributor: Unknown, 1 Screen",
  "Finding Dory, Distributors: Pixar and Disney, 8 Screens"
)
```

```
str_detect(lines, "Columbia|Pixar")
```

TRUE TRUE FALSE TRUE

Making things optional

```
str_view(lines, pattern = "Distributor|Distributors")
```

```
str_view(lines, pattern = "Distributors?")
```

Karate Kid 2, Distributor: Columbia, 58 Screens

Finding Nemo, Distributors: Pixar and Disney, 10 Screens

Finding Harmony, Distributor: Unknown, 1 Screen

Finding Dory, Distributors: Pixar and Disney, 8 Screens

Greedy vs. lazy

```
str_view("Toy Story 3 In Disney Digital 3D", ".*3")
```

.*3

Toy Story 3 In Disney Digital 3D

```
str_view("Toy Story 3 In Disney Digital 3D", ".*?3")
```

.*?3

Toy Story 3 In Disney Digital 3D

Let's practice!

INTERMEDIATE REGULAR EXPRESSIONS IN R