

CS 1520: Recitation 3

Regular Expressions

What is a Regex?

- A special text string for describing a search pattern.
- Applications
 - Did the user enter a proper email id, username, password, etc?
 - Anonymize all the phone numbers in a given dataset of users in Pittsburgh
- Cheat Sheet
 - <https://www.debuggex.com/cheatsheet/regex/javascript>
 - <https://www.cheatography.com/savagedesign/cheat-sheets/javascript-regexp/>

Examples

- Write a Regex for all the phone numbers in USA
 - (222) 345 4444, 232-446-4444, 212.566.7777, (676).677.8888,....
- `\(?:\d{3}\)?[-.]\d{3}[-.]\d{4}`
- But the above regex accepts phone numbers with non-matching parenthesis too, like (222 333 4444.
- Solution
 - `(\\d{3})[-.]\d{3}[-.]\d{4}|\\d{3}[-.]\d{3}[-.]\d{4}`

Examples

- Write a Regex for passwords that contain between 8-20 characters and can have, alphabets, special characters and digits
- `^[a-zA-Z -~]{8,20}$`

Dec	Hex	Oct	Chr	Dec	Hex	Oct	HTML	Chr	Dec	Hex	Oct	HTML	Chr	Dec	Hex	Oct	HTML	Chr
0	0	000	NULL	32	20	040	 	Space	64	40	100	@	@	96	60	140	`	`
1	1	001	SoH	33	21	041	!	!	65	41	101	A	A	97	61	141	a	a
2	2	002	SoTxt	34	22	042	"	"	66	42	102	B	B	98	62	142	b	b
3	3	003	EoTxt	35	23	043	#	#	67	43	103	C	C	99	63	143	c	c
4	4	004	EoT	36	24	044	$	\$	68	44	104	D	D	100	64	144	d	d
5	5	005	Enq	37	25	045	%	%	69	45	105	E	E	101	65	145	e	e
6	6	006	Ack	38	26	046	&	&	70	46	106	F	F	102	66	146	f	f
7	7	007	Bell	39	27	047	'	'	71	47	107	G	G	103	67	147	g	g
8	8	010	Bsp	40	28	050	((72	48	110	H	H	104	68	150	h	h
9	9	011	HTab	41	29	051))	73	49	111	I	I	105	69	151	i	i
10	A	012	LFeed	42	2A	052	*	*	74	4A	112	J	J	106	6A	152	j	j
11	B	013	VTAB	43	2B	053	+	+	75	4B	113	K	K	107	6B	153	k	k
12	C	014	FFeed	44	2C	054	,	,	76	4C	114	L	L	108	6C	154	l	l
13	D	015	CR	45	2D	055	-	-	77	4D	115	M	M	109	6D	155	m	m
14	E	016	SOut	46	2E	056	.	.	78	4E	116	N	N	110	6E	156	n	n
15	F	017	SIn	47	2F	057	/	/	79	4F	117	O	O	111	6F	157	o	o
16	10	020	DLE	48	30	060	0	0	80	50	120	P	P	112	70	160	p	p
17	11	021	DC1	49	31	061	1	1	81	51	121	Q	Q	113	71	161	q	q
18	12	022	DC2	50	32	062	2	2	82	52	122	R	R	114	72	162	r	r
19	13	023	DC3	51	33	063	3	3	83	53	123	S	S	115	73	163	s	s
20	14	024	DC4	52	34	064	4	4	84	54	124	T	T	116	74	164	t	t
21	15	025	NAck	53	35	065	5	5	85	55	125	U	U	117	75	165	u	u
22	16	026	Syn	54	36	066	6	6	86	56	126	V	V	118	76	166	v	v
23	17	027	EoTB	55	37	067	7	7	87	57	127	W	W	119	77	167	w	w
24	18	030	Can	56	38	070	8	8	88	58	130	X	X	120	78	170	x	x
25	19	031	EoM	57	39	071	9	9	89	59	131	Y	Y	121	79	171	y	y
26	1A	032	Sub	58	3A	072	:	:	90	5A	132	Z	Z	122	7A	172	z	z
27	1B	033	Esc	59	3B	073	;	;	91	5B	133	[[123	7B	173	{	{
28	1C	034	FSep	60	3C	074	<	<	92	5C	134	\	\	124	7C	174	|	
29	1D	035	GSep	61	3D	075	=	=	93	5D	135]]	125	7D	175	}	}
30	1E	036	RSep	62	3E	076	>	>	94	5E	136	^	^	126	7E	176	~	~
31	1F	037	USep	63	3F	077	?	?	95	5F	137	_	_	127	7F	177		Delete

charstable.com

Greedy and Lazy Evaluation

- Write Regex for all html tags in the given code snippet

`<div> This is CS 1520 </div>`

`<div> This is the last thing I am doing today </div>`

- Regex:

- `<.+>`

- `<.+?>`

Greedy and Lazy Evaluation

- Write Regex for all html tags in the given code snippet

`<div> This is CS 1520 </div>`

`<div> This is the last thing I am doing today </div>`

- Regex:

~~`<.+>`~~

- `<.+?>`

Regex For A Date

- Write the Regex for a date in the format yyyy/mm/dd or yyyy.mm.dd
- `\d{4}([\./])\d{2}([\./])\d{2}`
- Is this enough? Are all the matched dates valid?
 - No. Because it will also accept invalid dates like 1900/40/40

Regex For A Date

- Write the Regex for valid dates from 1900-01-01 through 2099-12-31
- `^(19|20)\d\d[- /.](0[1-9]|1[012])[- /.](0[1-9]|12)[0-9]|3[01])$`
- Is this enough?
 - No
 - Doesn't test for non-leap year dates, etc
 - Also, allows usage of inconsistent delimiters in the date

Regex Groups

Non-capturing group

- Final Validation of a date
- Capture the year, month and day and check for validation
- `^((?:19|20)\d\d)[- \\.](0[1-9]|1[012])[- \\.](0[1-9]|[12][0-9]|3[01])$`
- Here, we have
 - \$1 = year, \$2 = month, \$3 = day of the *yyyy-mm-dd* format date

Regex Groups

- Anonymize all the phone numbers in slide 3, in the format ddd-***-****
- `\(?:\d{3}\)?[-.](\d{3})[-.](\d{4})`
- Replace \$2 = '***' and \$3 = '****'

Backreferencing

- How to prevent dates with inconsistent delimiters?
 - Prevent dates like 1990/12.01
- `^(19|20)\d\d([- /.])(0[1-9]|1[012])\2(0[1-9]|12[0-9]|3[01])$`

Using Regex in JS – test()

- regex.test(str)

```
const string = "222-111-3344"  
const regex = /\(?\d{3}\)?[-. ]\d{3}[-. ]\d{4}/g;  
const doesExist = regex.test(string);  
console.log(doesExist);
```

- Output: true

Using Regex in JS – match()

- str.match(regex)
- Returns the matched strings

```
> var s = "my phone no are 111-2222 and 111-2233"
```

```
< undefined
```

```
> var r = /(\d{3})[-.]\d{4}/;
```

```
< undefined
```

```
> s.match(r)
```

```
< (2) ["111-2222", "111", index: 16, input: "my phone no are 111-2222 and 111-2233", groups: undefined]
```

```
> |
```

Problem with match()

- Does not return the groups with a global flag

```
> var s = "my phone no are 111-2222 and 111-2233"
```

```
< undefined
```

```
> var r = /(\d{3})[-.]\d{4}/g;
```

```
< undefined
```

```
> s.match(r)
```

```
< ▶ (2) ["111-2222", "111-2233"]
```

```
> |
```

Using Regex in JS – exec()

- regex.**exec**(str)
- Unlike match(), returns the groups even with a global flag

```
> var s = "my phone no are 111-2222 and 111-2233"
< undefined
> var r = /(\d{3})[.-]\d{4}/g;
< undefined
> r.exec(s)
< ▶ (2) ["111-2222", "111", index: 16, input: "my phone no are 111-2222 and 111-2233", groups: undefined]
> r.exec(s)
< ▶ (2) ["111-2233", "111", index: 29, input: "my phone no are 111-2222 and 111-2233", groups: undefined]
> r.exec(s)
< null
```


References

- <https://www.youtube.com/watch?v=7DG3kCDx53c&list=PLRqwX-V7Uu6YEypLuls7iidwHMdCM6o2w>
- <https://www.regular-expressions.info/dates.html>
- <https://www.talentcookie.com/2015/07/lets-practice-regular-expression/>
- <https://regex101.com/>