

# QSS20: Modern Statistical Computing

## Unit 08: Regular expressions (Regex)

# Goals for today

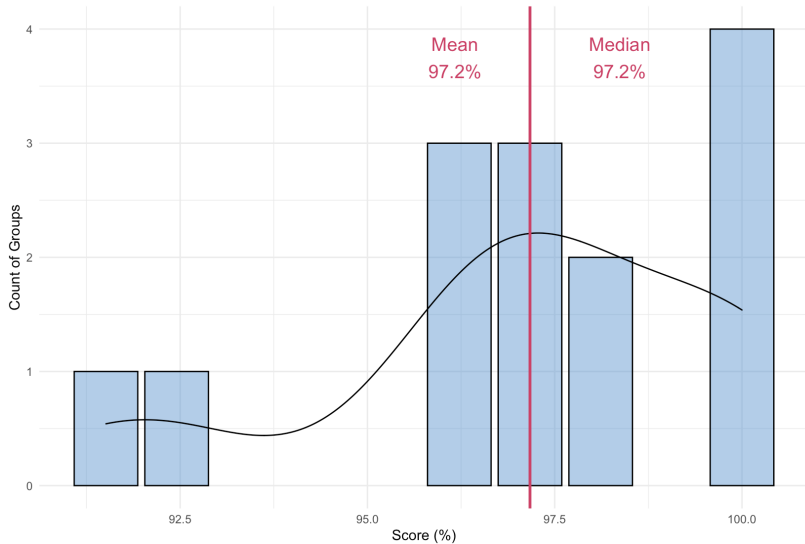
- ▶ Pset 2 feedback
- ▶ Recap of exact merging & LaTeX
- ▶ Regular expressions (Regex)

# Goals for today

- ▶ **Pset 2 feedback**
- ▶ Recap of exact merging & LaTeX
- ▶ Regex lecture & activity

# Pset 2 grade distribution

**Max:** 53.0; **Min:** 48.5

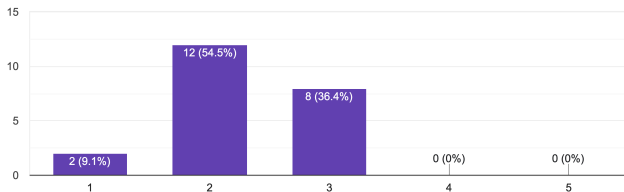


# Psets difficulty (1/5)

How would you rate the difficulty of problem set ONE (1)?

 Copy

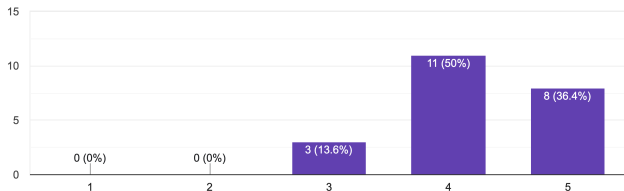
22 responses



How would you rate the difficulty of problem set TWO (2)?

 Copy

22 responses

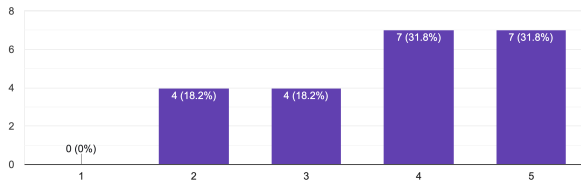


# Course learning vs. prep (2/5)

How much do you agree with this statement:  
"I am learning a lot from this course!"

 Copy

22 responses

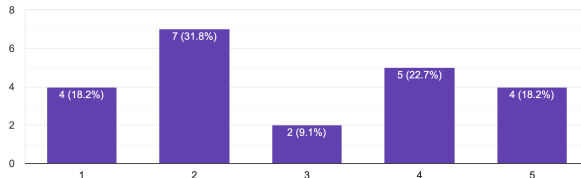


How much do you agree with this statement:

"I don't think I have the necessary background to do well in this course."

 Copy

22 responses

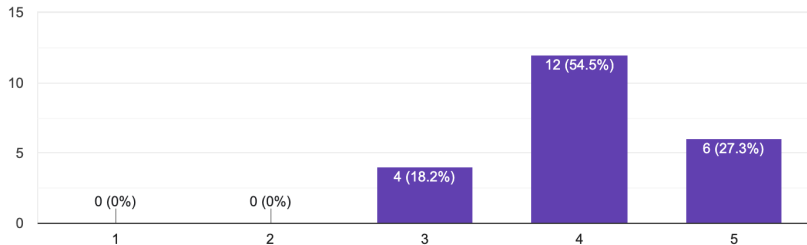


# Course pace (3/5)

How would you rate the overall pace of the course?



22 responses



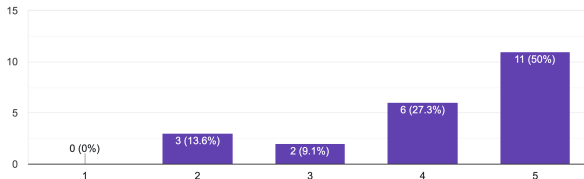
# Collaboration (4/5)

How much do you agree with this statement:

"Collaborating with a classmate resulted in a higher quality submission of pset 2."



22 responses

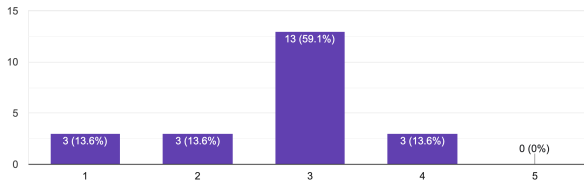


Between you and your classmate/partner, who did more work on your submission of pset 2?



(If it was balanced, choose the middle button, 3.)

22 responses





# Text comments & my response (5/5)

## **What's working:**

- ▶ Collaboration/camaraderie; partner psets
- ▶ Recaps
- ▶ Live coding

## **What to improve:**

- ▶ Too much time on recaps
- ▶ More guidance on in-class activities
- ▶ Take more time for lecture/code explanation
- ▶ Teach lambda functions
- ▶ Pset working sessions
- ▶ Tackle small problems one at a time; shorter in-class assignments

## **Changes to expect in class:**

- ▶ Quicker recaps
- ▶ More time on live coding, code explanation, working on code together
- ▶ X-Hours for pset work sessions
- ▶ Simpler, more modular in-class activities

# Join your final project group & pick a group name!

<b>Project</b>	<b>Partner A</b>	<b>Partner B</b>	<b>Partner C</b>	<b>Partner D</b>
<b>1: Felony sentencing</b>	Adin McAuliffe	Jordan Miller	Sam O'Donnell	Sam Winchester
<b>2: SIP: Medical IDD training</b>	Spencer Allen	Emma Elsbecker	Sabin Hart	Alex Ma
<b>3: SIP: SIRS</b>	John D'avanzo	Chloe Terestchenko	Ryan Wu	Mei Xu
<b>4: SIP: SIRS</b>	Leyla Jacoby	Keren Luo	Isabel Pantle	Jessie Wang
<b>5: Indep. project</b>	Alex Craig	Aaron Xie		
<b>6: Undergrad religiosity @DC</b>	Will Bryant			
<b>7: SIP: Medical training</b>	Aimen Abdulaziz	Esmeralda Abreu-Jerez	Bernardo Burnes Garza	Charles Knight
<b>8: SIP: Medical training</b>	Akshay Kelshiker	Edgar Ozu-zun	Shawn Yoon	Jeremy Rodriguez

# Where we are

- ▶ Pset 2 feedback
- ▶ **Recap of exact merging & LaTeX**
- ▶ Regular expressions (Regex)

# Recap of exact merging and LaTeX

What do you remember?

# Recap of exact merging

## Tips:

- ▶ Main data on left, aux on right
- ▶ Ideally, unique join key with same name in main/aux (can also use multiple cols)
- ▶ Four main kinds of joins: inner (shared keys), outer (all), left, right
- ▶ Check # rows before & after merge
- ▶ To debug, check for spelling variations in join key

## Useful commands:

```
pd.merge(maindf, auxdf, on='joinkey') # defaults to inner
pd.merge(maindf, auxdf, on='joinkey', how='outer') # drop no rows
pd.merge(maindf, auxdf, left_on='leftkey',
         right_on='rightkey') # diff keys
indicator = 'unit_mergesource' # merge diagnostic
suffixes = ('_main', '_aux') # for shared non-join cols
```

# Recap of LaTeX

## Tips:

1. Compile frequently
2. Be sure to exit any environments you made (e.g., `itemize`)
3. Mind the special characters:

escape character: backslash (`\`)

comment character: percent sign (`%`)

column separator: ampersand (`&`)

## Useful commands:

```
\begin{itemize} % start list
```

```
\begin{enumerate} % start NUMBERED list
```

```
\section{Title of section} \label{sec:shorthand}
```

```
\ref{sec:shorthand} % reference section in text
```

```
\begin{figure}
```

```
\includegraphics[scale = 0.5]{example_graphic.png}
```

*What's missing from these code snippets?*

# Where we are

- ▶ Pset 2 feedback
- ▶ Recap of exact merging & LaTeX
- ▶ **Regular expressions (Regex)**

# Today: basic regex to improve match rates for strings as join keys

- In example below, what if we didn't have the NCES ID numeric identifier? Ways to improve match rates for spelling variations (sometimes called `entity resolution`)

<b>Student</b>	<b>Year</b>	<b>District</b>
Jeremy	2021	New Trier High School
Emma	2022	Hanover High
Esmeralda	2022	Homeschool
⋮		

<b>District</b>	<b>% FRPL</b>
New Trier HS	X%
Hanover HS	Y%
Lebanon HS	Z%
⋮	



## Working example

Want to clean school names and classify them into different types (elementary; middle school; high school; charter; alternative; etc...). E.g.:

Central Columbia Ms

Riley County High

Jarrell H S

Trumbull School

SAN GABRIEL ELEMENTARY

Plains El

Pond Hill School

Franklin Elementary

P.S. 119

ANDREW CARNEGIE MIDDLE

Example of variety of names that all match the pattern within `str.contains`

```
1 cep_optin['is_elem'] =  
2 np.where(cep_optin.schoolname_lower.str.contains("\s+elem",  
3 regex = True), True, False)
```

Examples of True show a lot of variation that could make merges to other data difficult...

paint branch elementary

stewart county elementary school

stove prairie elementary school

winchester avenue elementary school

oak hill elem.

lewis and clark elem.

saunemin elem school

desert springs elementary school

fifth district elementary

linden elementary school

## re module provides more flexible alternative to pandas str methods

- ▶ Need to import at top: `import re`
- ▶ General structure: `re.something(r'‘somepattern’’, some_str)`
- ▶ Challenging part is constructing the pattern that captures what you want to capture

# Open activity notebook

*Work along with your final project group!*

Follow along the first part (before group activity) of [04\\_regex\\_blank.ipynb](#)

# First example

- We want to standardize the school names so that if it's an elementary school, it has the string “elemschool” after its other identifiers. E.g.:

original	cleaned
paint branch elementary	paint branch elemschool
stewart county elementary school	stewart county elemschool
stove prairie elementary school	:
winchester avenue elementary school	
oak hill elem.	
lewis and clark elem.	
saunemin elem school	
desert springs elementary school	
fifth district elementary	
linden elementary school	

# Approach 1: re.sub with stubborn listing

Basic syntax:

```
re.sub(pattern, str_to_sub_in, str_to_match)
```

**Step 1:** Construct a pattern to match all identified variations for which you want to substitute. Regex is very flexible, so there are *many ways to do this!* To start out, let's write out all the variations we can think of and stubbornly include them all with | (or), e.g.:

```
1 ## define pattern
2 elem_pattern = r"elementary|elem|elem\.|elementary school"
3
4 ## replace in one string
5 one_str_clean = re.sub(elem_pattern, "elemschool", one_str)
6
7 ## replace for all strings in the column schoolname_lower
8 all_str_clean = [re.sub(elem_pattern, "elemschool", one_str)
9                  for one_str in df.schoolname_lower]
```

But this exhaustive approach leads to issues...

# Limits to trying to include each option separately

(1) Sees elementary and subs it out but leaves school in; (2) leaves in . after elem. since matches elem

orig_name	cleaned_name
paint branch elementary	paint branch elemschool
stewart county elementary school	stewart county elemschool school
stove prairie elementary school	stove prairie elemschool school
winchester avenue elementary school	winchester avenue elemschool school
oak hill elem.	oak hill elemschool.
lewis and clark elem.	lewis and clark elemschool.
saunemin elem school	saunemin elemschool school
desert springs elementary school	desert springs elemschool school
fifth district elementary	fifth district elemschool
linden elementary school	linden elemschool school

# A better idea: Using “metacharacters” or shortcuts to match general types of patterns

## The basic metacharacters

Common ones:

- ▶ `\d` or `[0-9]`: numbers
- ▶ `[A-Z]`: uppercase alpha (any; can do subsets like `[ABC]`)
- ▶ `[a-z]`: lowercase alpha
- ▶ `\w`: alphanumeric (so numbers of letters)
- ▶ `+`: match one or more appearances (e.g., if we have several usernames—`jhaber-zz`; `jhaber-zzz`; `jhaber-zzzz`—could do `[a-z]+-z+` to match all versions)
- ▶ `*`: match any number
- ▶ `^`: match at beginning of string or line
- ▶ `$`: match at end of string or line
- ▶ `{x, y}`: match a pattern of length between `x` and `y` (e.g., to capture numbers that look like ages and could be length 1-3, write as `\d{1,3}`)



# Using metacharacters to make the previous pattern more robust to variations

```
1  ## define pattern
2  elem_pattern_try2 = r"(elem.*)"(\s+)?(school)?"
```

Breaking down each component:

- ▶ Creating groups using ( ): these help us define groups of characters to look at together
- ▶ elem.\*: matches elem, elem., and elementary (would maybe want to make more restrictive if we had schools named things like element that we didn't want to match)
- ▶ Multiple spaces: uses “metacharacter” to match 1+ spaces:  
 \s+
- ▶ ? : **optional pattern**, or match if this pattern occurs but also match if it doesn't (in this case, it's useful since schools like fifth district elementary have nothing afterwards)
- ▶ (school)? : similarly, sometimes ends with school and we want to replace; other times just ends with elementary or elem

# Executing re.sub

```
1 ## define pattern
2 elem_pattern_try2 = r"(elem.*)(\s+)?(school)?"
3
4 ## replace for all strings in the column schoolname_lower
5 all_str_clean_try2 = [re.sub(elem_pattern_try2, "elemschool",
6                             one_str) for one_str in df.schoolname_lower]
7
8 ## then assign the resulting list to the df
9 df['cleaned_name_try2'] = all_str_clean_try2
10
11 ## could easily use same approach to create new column directly
12 df['cleaned_name_try2'] = df['schoolname_lower'].apply(
13     lambda name:
14         re.sub(elem_pattern_try2, "elemschool", name))
```

How does the result compare to our first, more exhaustive method?

<b>orig_name</b>	<b>cleaned_name</b>
paint branch elementary	paint branch elemschool
stewart county elementary school	stewart county elemschool
stove prairie elementary school	stove prairie elemschool
winchester avenue elementary school	winchester avenue elemschool
oak hill elem.	oak hill elemschool
lewis and clark elem.	lewis and clark elemschool
saunemin elem school	saunemin elemschool
desert springs elementary school	desert springs elemschool
fifth district elementary	fifth district elemschool
linden elementary school	linden elemschool

# Other re operations

- ▶ In previous example, we:
  1. Defined a pattern: different variations of the string 'elementary'
  2. Used `re.sub(pattern, replacement, string)` to substitute the target string with something else (in this case, 'elemschool') when we match the pattern
- ▶ In other examples, we may want to:
  - ▶ Define a pattern that characterizes different subparts of a string (e.g., maybe 'elementary' is one part we want to extract, and 'school' is another)
  - ▶ Match that pattern
  - ▶ Extract the matches and do something

# Three similar matching methods for regex

## 1. `re.findall(pattern, string)`

- ▶ **What it does/returns:** scans the string from left to right and returns the matches in a **list**
- ▶ **Useful for:** parsing a string and then recombining elements (e.g., elementary could be list element 0; school could be list element 1)

## 2. `re.search(pattern, string)`

- ▶ **What it does/returns:** scans the **entire string** and returns the substring(s) regardless of where it appears in the string. If no matches found, returns `None`. If matches found, returns **MatchObject**
- ▶ **Useful for:** checking *if* there's a match (since can check whether result is equal to `None`); same use above of getting substrings

## 3. `re.match(pattern, string)`

- ▶ **What it does/returns:** Operates similarly to `re.search()` but rather than searching the entire string, only returns if found at beginning of string

Both `re.match()` and `re.search()` only return the first appearance of a pattern; if want to return all occurrences (e.g., count how many times “Trump” appears in a single tweet), can use either `re.findall()` or `re.finditer()`

# Executing re.findall

**Example:** match words before and after the word 'charter'

```
1 ## define charter pattern
2 charter_pattern = r"(.*)(\s+(charter)(\s+)?(\w+)?"
3
4 ## execute re.findall
5 test_charter_findall = [re.findall(charter_pattern, school)
6                          for school in charter_ex.schoolname]
```

# What this returns: a list of lists

## orig\_name

buffalo collegiate charter school  
 thomas edison charter academy  
 moving everest charter school  
 life source international charter  
 south valley academy charter school  
 neighborhood charter school of harle  
 brighter choice charter school-girls  
 children's community charter  
 frontier elementary school  
 columbus humanities, arts and...  
 okemos public montessori-central  
 pawhuska es  
 east valley senior high  
 glenpool es  
 number 27  
 south fork elementary

```

: [[('buffalo collegiate', 'charter', ' ', 'school')],
  [('thomas edison', 'charter', ' ', 'academy')],
  [('moving everest', 'charter', ' ', 'school')],
  [('life source international', 'charter', ' ', 'school')],
  [('south valley academy', 'charter', ' ', 'school')],
  [('neighborhood', 'charter', ' ', 'school')],
  [('brighter choice', 'charter', ' ', 'school')],
  [('children's community', 'charter', ' ', 'school')],
  [],
  [],
  [],
  [],
  [],
  [],
  [],
  [],
  []]
```

## Executing `re.search` using same pattern and school names list

```
1 ## charter pattern
2 charter_pattern = r"(.*)\s+(charter)(\s+)?(\w+)?"
3
4 ## search
5 test_charter_search = [re.search(charter_pattern, school)
6                         for school in charter_ex.schoolname]
```



`re.search` returns `MatchObject` if match found, else `None`

### **orig\_name**

buffalo collegiate charter school

thomas edison charter academy

moving everest charter school

life source international charter

south valley academy charter school

neighborhood charter school of harle

brighter choice charter school-girls

children's community charter

frontier elementary school

columbus humanities, arts and...

okemos public montessori-central

pawhuska es

east valley senior high

glenpool es

number 27

south fork elementary

```
[<re.Match object; span=(0, 33), match='buffalo collegiate
<re.Match object; span=(0, 29), match='thomas edison char
<re.Match object; span=(0, 29), match='moving everest cha
<re.Match object; span=(0, 33), match='life source intern
<re.Match object; span=(0, 35), match='south valley acade
<re.Match object; span=(0, 27), match='neighborhood char
<re.Match object; span=(0, 30), match='brighter choice ch
<re.Match object; span=(0, 28), match="children's communi
None,
None,
None,
None,
None,
None,
None,
None,
None]
```

# Extracting matches from MatchObject with .group() method

Example: thomas edison charter academy

```
1 ##### here , we're just focusing on the 3rd match = 6th entry (thomas
   edison charter academy)
2 ##### and we're getting the first group from that match
3 thomas_match = test_charter_search[5]
4 thomas_match
5
6 ##### example where we're just getting the first group
7 ##### (name of school before charter)
8 thomas_firstgroup = thomas_match.group(1)
```

Prints: thomas edison

## Show all groups for that one match

```
1 ### iterate over all groups and print
2 for i in range(0, len(thomas_match.groups())+1):
3     print("Group " + str(i) + " is: ")
4     print(thomas_match.group(i))
```

Prints:

```
1 Group 0 is:
2 thomas edison charter academy
3 Group 1 is:
4 thomas edison
5 Group 2 is:
6 charter
7 Group 3 is:
8
9 Group 4 is:
10 academy
```

Can also extract the groups as a tuple:

```
1 thomas_match.groups()

1 ('thomas edison', 'charter', ' ', 'academy')
```

## Break for activity

*Collaborate with your final project group!*

**Group activity** section at end of [04\\_regex\\_blank.ipynb](#)