## Congratulations! You passed!

blue (?=whale)

**⊘** Correct

Grade received 28.57% To pass 10% or higher

Go to next item

Reflecting back on what you learned Total points 14 1. Hi all! 0 / 1 point  $Welcome\ to\ SIADS\ 521.\ Let's\ start\ the\ course\ by\ reviewing\ what\ you\ learned\ in\ SIADS\ 505.\ This\ will\ help\ you\ with$  $hand ling\ data\ for\ visualizations\ you\ will\ create\ throughout\ this\ course.\ This\ will\ also\ help\ me\ understand\ your$ background knowledge. This will not affect your grade! This is just so I get a sense as to how sticky the items you learned in SIADS 505 are! Now, let's dig in! Which regular expression will return all 'A's and 'B's from the string below? grades="AAAABBBCCCCCDDDD" O [AB] [A][B] (A&B) (AB) igotimes Incorrect This will capture A followed by B.  $\textbf{2.} \quad \text{In a regular expression, what is the special character that matches zero or more characters?}$ 0 / 1 point • + O \* O ? O # 0 ^ + captures one or more characters. 3. Which regular expression pattern would find the word 'blue' only when it is followed by 'whale' and NOT when 1/1 point followed by 'bottle' in the sentence "The blue whale ate the blue bottle."? O blue .whale blue (.whale) blue (?P<whale>) blue [whale]

. What	will be re	eturne	d after y	ou run	the regu	ılar exp	oressio	on below?							0/1 point
resul	ts = re.ma	atch(r'(	\w+)@(\	w+).(\v	/+)', 'user	id@um	nich.ed	lu')							
print	results.g	roups(	))												
<b>(1)</b>	userid@u	ımich.e	<u>edu</u> ′												
0 1	umich'														
0 1	userid'														
0 1	edu'														
('userid', 'umich', 'edu')															
Incorrect  To have this pattern with re.match, you can use '(\w+@\w+.\w+)'															
	h statem			bout th	ne code k	oelow?									0 / 1 point
#	port pa														
		ataFr			: ['31/ ': [2,			1/9/2021', '	2/9/	2021']	,				
	<pre># code2 df['date'] = pd.to_datetime(df['date'], dayfirst=True)</pre>														
	After you run code1 and sort the dataframe df based on the 'date', the sequence of the values in the items column will be still 2, 3, and 4														
O F	After you run code1 only, the data type of the values in the date column is the integer.														
O F	After you run code1 and code2, the data type of the values in the date columns is the object.														
		code2	with df	['date']	= pd.to_	datetin	ne(df['c	date'], dayfirst=l	alse)	would	not m	ake a (	differenc	e in	
	output. The outpu	ıt of df	chano	c (2 2)											
	The outpo	at or ur	.зпаре	3 (2, 3)											
$\otimes$	Incorrect If dayfire 9th of 2	st=Fals	e, then	the "1/	9/2021" a	and "2/	9/2021	" will be recogn	ized a	as Janua	ary 9th	of 20	21, and	February	,
. Whic	h of the f	ollowin	ng optio	ns will	merge df	1 and	df2 on	the left into df_	merg	ed on t	he rigi	nt?			0 / 1 point
	Key	Α	В		Key	С	D			Key	Α	В	С	D	
0	K0	A0	В0	0	K0	C4	D4		0	K0	A0	В0	C4	D4	
1	K1	A1	B1	1	K1	C5	D5		1	K1	A1	В1	C5	D5	
2	K2	A2	B2	2	K4	C6	D6		2	K2 K3	A2 A3	B2	NaN NaN	NaN NaN	
3		A3	ВЗ	3	K5	C7	D7		3	No		ВЗ		INGIN	
	dí	1			df	2					df_	merc	jed		
pd.merge(df1, df2, on="Key", how="outer")															
● F	od.merge	(df1, d	f2, on="	Key", h	ow="inn	er")									
O F	od.merge	(df1, d	f2, on="	Key", h	ow="righ	nt")									
O F	od.merge	(df1, d	f2, on="	Key", h	ow="left"	")									
pd.merge(df1, df2, on="Key", how="center")															

(Note that the state of the sta

7.

Look at the code and the first DataFrame output about cereals below. Which code should you insert to get the second DataFrame output showing means and sums of the calories per manufacturer which is represented by the column "mfr"?

```
import pandas as pd
df = pd.read_csv('cereal.csv')
df.head()
```

	name	mfr	type	calories	protein	fat	sodium	fiber	carbo	sugars	potass	vitamins	shelf
0	100% Bran	N	С	70	4	1	130	10.0	5.0	6	280	25	3
1	100% Natural Bran	Q	С	120	3	5	15	2.0	8.0	8	135	0	3
2	All- Bran	K	С	70	4	1	260	9.0	7.0	5	320	25	3
3	All- Bran with Extra Fiber	К	С	50	4	0	140	14.0	8.0	0	330	25	3
4	Almond Delight	R	С	110	2	2	200	1.0	14.0	8	-1	25	3

## ### INSERT CODE

## mean sum

mfr		
Α	100.000000	100
G	111.363636	2450
K	108.695652	2500
N	86.666667	520
Р	108.888889	980
Q	95.000000	760
R	115.000000	920

- df.groupby("calories")["mfr"].aggregate(np.mean, np.sum)
- $\\ \bigcirc \ \ df.groupby("calories")["mfr"].agg([lambda \ x : x.mean(), lambda \ x : x.sum()]) \\$
- $\begin{tabular}{ll} \hline & df.groupby("mfr").aggregate({"calories":np.mean, np.sum}) \\ \hline \end{tabular}$
- df.groupby("mfr").agg({"calories":[np.mean, np.sum]})
- df.groupby("mfr").transform({"calories":[np.mean, np.sum]})
  - ✓ Correc

df.groupby(the column you want to use for the groupby).agg({the column you want to apply np.mean and np.sum: [np.mean, np.sum]})

8. What is the right output of the following code?

1/1 point

```
        col1
        col2

        0
        False
        True

        1
        False
        False

        2
        True
        False

        3
        False
        True

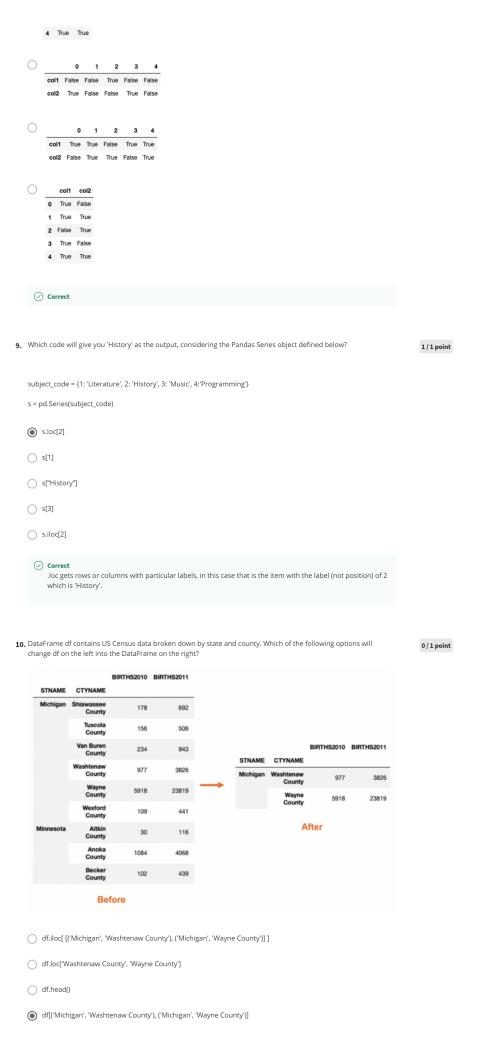
        4
        False
        False
```

col1 col2

```
    False True
    True
```

2 True False

3 False True



⊘ df.loc		ashtenaw County), ('Michigan', 'Wayne County')]	
The	code will thro	w an error since ('Michigan', 'Washtenaw County') and ('Michigan', 'Wayne County') will be keys to look for.	
11. The follow	ving figure sho	ws the Series s. If you run the code below, what will be its output?	0 / 1 poi
0 0.1 1 1.2 2 2.3 3 3.4 4 Né 5 Né 6 6,7	. 0 . 0 . 0 an an . 0		
	float64		
s = s.fillna s	(method="ffill")		
© 0 1 2 3 4 5 6 7 dty	0.0 1.0 2.0 3.0 4.0 5.0 6.0 7.0 pe: float64		
0 1 2 3 4 5 6 6 7 dt	0.0 1.0 2.0 3.0 3.0 3.0 6.0 7.0 ype: float6	4	
0 1 2 3 4 5 6 7 dty	0.0 1.0 2.0 3.0 6.0 6.0 7.0 pe: float64		
0 1 2 3 4 5 6 7 dtys	0.0 1.0 2.0 3.0 NaN NaN 6.0 7.0		
0 1 2 3 4 5 6 7 dtyp	0.0 1.0 2.0 3.0 3.0 6.0 6.0 7.0 e: float64		
⊗ Inco		l") will fill out NaN values following values of the previous rows.	
12 Which of	he following o	ptions is a correct statement about the DataFrame df?	
		Score	0 / 1 poi
school1	Heru History	86	
school3	Rene Physics	92	
school2	Alex Biology Yuna Math	76 72	
30110012	· · · · · · · · · · · · · · · · · · ·		

O You can have only one column as an index.

- (Name') will drop the column 'Name'. df.loc[:, ["Name", "Score"]] will give you all Names and Scores for all schools. O You cannot set mixed types of data values as index labels. igotimes Incorrect The default setting of . drop method is dropping a row. If you would like to drop a column, you need to setthe parameter "axis=1" 13. Which of the following options is a correct statement about the code below? : import re test\_string = "ABCDBCBCBCBBBBDE"
  pattern = r"pattern" re.findall(pattern, test\_string) None above is correct. Including a non-greedy quantifier in the pattern is one way to get the output above.

  - $\bullet$  With the pattern = r"(B)\*", you can get the same output.
  - If re.findall(pattern, test\_string) is replaced with re.finditer(pattern, test\_string), you cannot get an output in a
  - $\bigcirc$  With the pattern = r"B+", you can get the output above.
  - igotimes Incorrect

With  ${\rm *}$  quantifier which is for zero or more, you will get spaces too for matches.

14. The df on the left side below was generated after running code 1. Then, after running code 2, the df has changed as  $the \ df \ on \ the \ right \ side \ below. \ With \ that \ said, \ which \ of \ the \ following \ options \ is \ an \ INCORRECT \ statement \ about \ the$ code and the df below?

0 / 1 point

0 / 1 point

```
import pandas as pd
### code1
df.rolling(2).mean()
```

	Temp		Tem
Date		Date	
1981-01-01	20.7	1981-01-01	NaN
1981-01-02	17.9	1981-01-02	19.30
1981-01-03	18.8	1981-01-03	18.35
1981-01-04	14.6	1981-01-04	16.70
1981-01-05	15.8	After running code 2	15.20
1990-12-27	14.0	1990-12-27	14.30
1990-12-28	13.6	1990-12-28	13.80
1990-12-29	13.5	1990-12-29	13.55
1990-12-30	15.7	1990-12-30	14.60
1990-12-31	13.0	1990-12-31	14.35
3650 rows ×	1 colu	mns 3650 rows >	1 col

- O None of the above is incorrect.
- $\bigcirc \ \ \, \text{The reason why the first value of the df on the right side is NaN is because the rolling window was larger than}$ 1 and therefore needed more values to calculate the mean.
- To calculate means for every two days as above, we can also apply resampling() to get the similar result.
- df.rolling(2).mean() could be written as df.rolling("2D").mean() since the values in the Date column are parsed as the datetime data type.
- The value 2 in the rolling(2) defines the size of rolling windows, in this case, which is every 2 rows.