

WRITTEN TEST for DATA QUALITY ANALYST @ TURBOLAB TECHNOLOGIES

Swaroop N C

E-mail : ncswaroop1997@gmail.com

Mobile : +91 99 1616 0739

LinkedIn : <https://www.linkedin.com/in/swaroop-n-c-a09499131/>

Date : 09-05-2022

Instructions: 3 hours to complete

- All three questions should be attempted.

```
1 import numpy as np
2 import pandas as pd
```

▼ QUESTION 1

Part 01. **Create a Pandas series with 100 random dates as it falls between 01-01-20 to 01-01-21**

(hint: use **Pandas date_range** function).

Like:- _ 2020-01-06 2020-06-11 2020-02-18

Part 02. **Dedupe it and calculate the number of duplicates and convert it to percentage.**

Then by **using regex**, filter values where either the month is 02,05,09 OR the date is 01,04,07 - the apply function should not be used.

Finally, **calculate the percentage of values filtered for the month condition and the date condition.**

▼ SOLUTION TO Q 1

```
1 import random
```

```
1 start_date='1/1/2020'
2 end_date='1/01/2021'
3 n=100
```

```
1 All_dates=pd.date_range(start=start_date, end=end_date)
```

2 All_dates

```
DatetimeIndex(['2020-01-01', '2020-01-02', '2020-01-03', '2020-01-04',
               '2020-01-05', '2020-01-06', '2020-01-07', '2020-01-08',
               '2020-01-09', '2020-01-10',
               ...,
               '2020-12-23', '2020-12-24', '2020-12-25', '2020-12-26',
               '2020-12-27', '2020-12-28', '2020-12-29', '2020-12-30',
               '2020-12-31', '2021-01-01'],
              dtype='datetime64[ns]', length=367, freq='D')
```

```
1 df= pd.DataFrame(random.choices(x,k=n))
2 df
```

☞ **0** 

0 2020-01-23

1 2020-08-25

2 2020-05-12

3 2020-04-18

4 2020-07-24

...

95 2020-03-17

96 2020-09-18

97 2020-10-04

98 2020-05-01

99 2020-09-06

100 rows × 1 columns

```
1 df.columns=['Date']
```

```
1 df
```

	Date
0	2020-01-23
1	2020-08-25
2	2020-05-12
3	2020-04-18



Dedupe it and calculate the number of duplicates and convert it to percentage.

... ...

```
1 df.nunique()
```

```
Date      88
dtype: int64
```

```
1 Num_Of_Duplicates = n - df.nunique()
2 Num_Of_Duplicates
```

```
Date      12
dtype: int64
```

```
1 Dup_Percent = ((Num_Of_Duplicates)/n)*100
2 Dup_Percent
```

```
Date      12.0
dtype: float64
```

Then by using regex, filter values where either the month is 02,05,09 OR the date is 01,04,07 - the apply function should not be used.

Finally, calculate the percentage of values filtered for the month condition and the date condition.

```
1 #COULD NOT ANSWER THIS PART
```

▼ QUESTION 2

This question has two parts,

a) Sentence Validator and

b) Name Reducer

a) Create a sentence validator function. The validator should return input if it is valid otherwise return False

Validation Criteria:

1. Start letter must be an uppercase letter and it should follow either a lowercase letter or a single whitespace.
2. All letters in the sentence except the start letter must be in lowercase.
3. The last character (aka terminal character) of the sentence must be any of the following:
. (dot) ? (question mark) ! (exclamation mark)
4. Words must be separated with a single whitespace.

If there is a hyphen between any two words then there should be one whitespace before and after that hyphen.

eg: Lab - 1 is valid, but Lab- 7 and Lab -7 are invalid

b) Write a reducer function to clean the sentences. Reducer takes output from the validator function as input and performs the following cleaning steps,

- Removes terminal characters (see above validation criteria for list of allowed terminal characters)
- Removes all duplicated word groups (see below examples) but keep its first occurrence
- Removes all leading and trailing whitespaces and hyphens

After completing the functions for part a & b.

Create a function `check_and_clean` which takes a sentence as input and validates (using validator function) and returns reduced string (output from reducer function) when the sentence is valid, "" otherwise.

Note: You are not allowed to use regex for this question.

Example 1: Input: Melo diagnostics melo Labs

Sentence is invalid, failed validation criteria 2 & 3

Output: invalid

Example 2:

Input: Melo diagnostics - southpark east 29th street - southpark east 29th street - free drug testing not offered. Sentence is valid

Word groups: Melo diagnostics southpark east 29th street free drug testing not offered

Output: Melo diagnostics - southpark east 29th street - free drug testing not offered

Example 3:

Input: Simple labs - covid test available - west side hospital - covid test available!

Sentence is valid

Word groups: Simple labs covid test available west side hospital

Output: Simple labs - covid test available - west side hospital

▼ SOLUTION TO Q 2

sentence validator

Validation Criteria:

1 Start letter must be an uppercase letter and it should follow either a lowercase letter or a single whitespace.

2 All letters in the sentence except the start letter must be in lowercase.

3 The last character (aka terminal character) of the sentence must be any of the following: . (dot) ? (question mark) ! (exclamation mark)

Words must be separated with a single whitespace. If there is a hyphen between any two words then there should be one whitespace before and after that hyphen.

eg: Lab - 1 is valid, but Lab- 7 and Lab -7 are invalid

```

1 def SentenceValidator(string):
2     length = len(string)
3     #Start letter must be an uppercase letter
4     if (string[0].islower()):
5         return False
6     #last character of sentence must be . (dot) ? (question mark) ! (exclamation mark)
7     elif not(string[length-1] == '.' or string[length-1] == '?' or string[length-1] == '!
8         return False
9     #All letters in the sentence except the start letter must be in lowercase
10    else:
11        for ele in string[1:]:
12            if ele.isupper():
13                return False
14    #Words must be separated with a single whitespace.
15    #check for more than 2 consecutive whitespace
16    for i in range(len(string)) :
17        if (string[i]==' ' and string[i+1]==' ') :
18            return False
19    #If there is a hyphen between any two words then
20    #there should be one whitespace before and after that hyphen.
21    for i in range(len(string)) :
22        if (string[i]=='-') :
23            if not((string[i-1]==' ')and (string[i+1]==' ')):
24                return False
25    return string

```

b) Write a reducer function to clean the sentences. Reducer takes output from the validator function as input and performs the following cleaning steps,

- **Removes terminal characters (see above validation criteria for list of allowed terminal characters)**
- **Removes all duplicated word groups (see below examples) but keep its first occurrence**

```

1 def reducer(ValidOp) :
2     #Removes terminal characters
3     ValidOp=ValidOp[:-1]
4     #Removes all leading and trailing whitespaces
5     ValidOp=ValidOp.strip()
6     #Removes all leading and trailing HYPHENS
7     ValidOp=ValidOp.strip("-")
8     #Removes all duplicated word groups but keep its first occurrence
9     l = ValidOp.split()
10    k = []
11    for i in l:
12        if (ValidOp.count(i)>=1 and (i not in k) or i =='-'):
13            k.append(i)
14    ValidOp=' '.join(k)
15    return ValidOp

```

check_and_clean

Create a function check_and_clean which takes a sentence as input and validates (using validator function) and returns reduced string (output from reducer function) when the sentence is valid, "" otherwise.

```

1 def check_and_clean() :
2     stringcheck=str(input("Input String : "))
3     if (SentenceValidator(stringcheck) != False ) :
4         print("***** OUTPUT *****")
5         return (reducer(SentenceValidator(stringcheck)))
6     else :
7         print("***** OUTPUT *****")
8         return ("<invalid>")

```

```

1 #driverprogram
2 #test case 1 Melo diagnostics melo Labs
3 check_and_clean()

```

```

Input String : Melo diagnostics melo Labs
***** OUTPUT *****
'<invalid>'

```

```

1 #driverprogram
2 #test case 3 : Simple labs - covid test available - west side hospital - covid test av
3 check_and_clean()

```

```
Input String : Simple labs - covid test available - west side hospital - covid test a
***** OUTPUT *****
'Simple labs - covid test available - west side hospital - '
```

▼ Question 3:

A sample data of posts of random users are given in this link: [click here to download](#) Post_id, date of post and post caption details are available in the sample dataset.

Create a function that extracts posts older than 13/11/2021 and finds the 3 most frequently used special characters out of it. The function should return the 3 most frequently used special characters and the number of times they occurred in the filtered data.

Example

post_id	date	caption
post #1	11/11/2021	@bla bl@ bla! 23 🔥
post #2	15/11/2021	Foo b@r foOB!a 🤖
post #3	12/11/2021	🔥 aerrt!! Qwe r rr\$
post #4	13/11/2021	@momo bati\$t@ 🔥

Output

```
[ (@, 4), (!, 3), (🔥, 2) ]
```

Explanation The post #2 is eliminated since it is older than the date 13/11/2021. In the remaining 3 rows, the special character “@” occurred the most i.e, 4 times in the posts #1 and #4. The second most frequent special character is “!” which occurred 3 times and then “🔥” occurred twice. Since only the top 3 most frequent ones are required, the remaining special characters “🤖” and “\$” are ignored.

```
1 #load dataset
2 df = pd.read_csv("Captions.csv")
3 df
```

	post_id	date	caption	
0	post #1	15/11/2021		NaN
1	post #2	14/11/2021	Skippers at the ground for a photo call but ca...	
2	post #3	18/11/2021	Kia ora everyone, We tried but it wasn't to be...	
3	post #4	11/11/2021	Teihorangi Walden 🏡	2022
4	post #5	13/11/2021		NaN
...
95	post #96	17/11/2021	Pink maomao, blue maomao, granddaddy hāpuku (N...	
96	post #97	12/11/2021	Hiking in Queenstown? Yes, please! 🥾	Now th...
97	post #98	11/11/2021	Opportunity to engage: Gender and Sex Diverse ...	
98	post #99	11/11/2021	7 years signed, sealed and delivered	❤️💙

extracting posts older than 13/11/2021

100 rows x 5 columns

```
1 df=df[df['date']<'13/11/2021']
2 df
```


	post_id	date	caption	
3	post #4	11/11/2021	Teihorangi Walden 🏠 2022	
5	post #6	12/11/2021	Lineup 🔥	
6	post #7	11/11/2021	"Taking some Mospiration from the hedges" @sam...	
8	post #9	12/11/2021	A reduction in wild animals, vegetation recove...	
10	post #11	12/11/2021	Which witch will you watch first this weekend?...	
13	post #14	11/11/2021	Just returned from the best kingfish fishing w...	
14	post #15	11/11/2021	ICYMI Our short film 'The Mountain Storm' wa...	
16	post #17	12/11/2021	The PENN BATTLE III DX Reel has landed! Exclus...	
18	post #19	11/11/2021	Excited to be back @nrl_weststigers as an offi...	
20	post #21	12/11/2021	🔥🔪🏠🏆 It all comes down to this! A trans-Tasman...	
22	post #23	12/11/2021	Working on our fitness influencer poses. Swipe...	
23	post #24	11/11/2021	🐘 vs 🐘🏠 Team Milestones: 🧡 Stephen Perofeta w...	
28	post #29	12/11/2021	NZKS is open for shopping this weekend! If you...	
31	post #32	12/11/2021	Twins Otukinekina and Valingi Kepu will join o...	
38	post #39	11/11/2021	Need warmth and performance? All the natural b...	
43	post #44	12/11/2021	👉👉👉👉	
46	post #47	12/11/2021	Grab a pack while the stock lasts! Get in quic...	
47	post #48	11/11/2021	Congratulations to the @volcanichillswinery Te...	
51	post #52	11/11/2021	2 years ago today we launched Live Ocean. Afte...	
52	post #53	12/11/2021	Grateful 🧡🧡	

```

1 def special_char_list(string):
2     special_char = 0
3     spec_char_list = []
4     for i in range(len(string)):
5         if(string[i].isalpha()):
6             continue
7         elif(string[i].isdigit()) :
8             continue
9         else :
10            if (string[i] not in spec_char_list) :
11                spec_char_list.append(string[i])
12    return(spec_char_list)

```

76 post #77 11/11/2021 Summer holidays on your mind? 🌞 we've got you ...

```

1 caption_text = ''
2 for row in df.caption :
3     caption_text += row

```

83 post #84 12/11/2021 🏒LIVE HOCKEY TOMORROW 🏒 - Link in our bio for l...

```

1 caption_text

```

▲

▼

▼

▼

▲

```

+ ',
☀️ ',
'E',
👤 ',
',
'\u200d',
♂',
🔪 ',
♥️ ',
📦 ',
',
👢 ',
'\u2060',
👉 ',
↑',
💙']

```

```

1 i=0
2 count_num=0
3 count=[]
4 for sp_char in spec_char_list :
5     for char in caption_text :
6         if(sp_char == char) :
7             count_num+=1
8     count.append(count_num)
9     count_num=0

```

```
1 len(count)
```

79

```
1 len(spec_char_list)
```

79

```

1 zip_ob = list(zip(spec_char_list,count))
2 zip_ob

```

```

(':', 10),
('/', 7),
('📷', 1),
('?', 7),
('💩', 1),
('&', 6),
('\u2063', 7),
('||', 2),
('"""', 10),
('»', 5),
('$', 2),
('💣', 1),
('🔪', 1),
('👤', 3),
('🏆', 2),
('*', 1),
('📷', 5),
('\'', 23),
('👉', 1)

```

```
( '👉', 1),
( '👈', 1),
( '❤️', 2),
( '🚗', 1),
( '...', 5),
( '😞', 1),
( '🐦', 1),
( '🌊', 5),
( '🔹', 2),
( '🔹', 2),
( '\U0001f90d', 1),
( '💪', 1),
( '‘', 1),
( '👉', 2),
( ' ', 1),
( '—', 1),
( '🍦', 1),
( 'I', 2),
( 'N', 2),
( 'Z', 1),
( '👉', 2),
( '🌀', 1),
( '😊', 1),
( '❤️', 4),
( '+', 1),
( '☀️', 1),
( 'E', 1),
( '👤', 1),
( ' ', 1),
( '\u200d', 1),
( '♂', 1),
( '🔪', 2),
( '❤️', 1),
( '📦', 1),
( ' ', 1),
( '👤', 1),
( '\u2060', 4),
( '👉', 1),
```

```
( '↑', 1),
( '...', 1),
```

1 sorted(zip_ob)

```
( '...', 5),
( '»', 5),
( '\u200d', 1),
( '—', 5),
( '—', 1),
( '‘', 1),
( '’', 23),
( '“', 2),
( '”', 2),
( '\u2060', 4),
( '\u2063', 7),
( '♂', 1),
( '❤️', 4),
( '—', 1),
( '↑', 1),
( ' ', 7),
( 'E', 1),
( '...', 1),
```

```
( 'I', 2),
( 'N', 2),
( 'Z', 1),
( '🌊', 5),
( '☀️', 1),
( '🌿', 1),
( '🍌', 1),
( '🏠', 1),
( '🏆', 2),
( '🔧', 1),
( '🔪', 2),
( ' ', 1),
( ' ', 1),
( ' ', 1),
( '🐘', 1),
( '🦎', 1),
( '🐭', 1),
( '👉', 1),
( '👊', 2),
( '👊', 1),
( '💙', 1),
( '💚', 1),
( '💛', 2),
( '💣', 1),
( '💪', 1),
( '📷', 1),
( '📷', 5),
( '🔒', 1),
( '👁️', 1),
( '🔥', 1),
( '💎', 2),
( '💎', 2),
( '❤️', 7),
( '😬', 1),
( '👤', 1),
( '👥', 3),
( '🚗', 1),
( '\U0001f90d', 1),
( '😬', 1),
( '👉', 2),
( '👉', 1),
( '👉', 1)
```

```
1 res = list(reversed(sorted(zip_ob, key = lambda x: x[1])))
2 res
```

```
( '📷', 5),
( '»', 5),
( '—', 5),
( '\u2060', 4),
( '❤️', 4),
( '👥', 3),
( '🔪', 2),
( '👊', 2),
( 'N', 2),
( 'I', 2),
( '👉', 2),
( '💎', 2),
( '💎', 2),
( '💛', 2),
( '👉', 2)
```

```
( '🕒', 2),
( '$', 2),
( '|', 2),
( '"', 2),
( '"""', 2),
( '"', 2),
( '❤️', 1),
( '↑', 1),
( '👊', 1),
( '👢', 1),
( ' ', 1),
( '📺', 1),
( '❤️', 1),
( '♂', 1),
( '\u200d', 1),
( ' ', 1),
( '👤', 1),
( 'E', 1),
( '☀️', 1),
( '+', 1),
( '😬', 1),
( '👁️', 1),
( 'z', 1),
( '🎉', 1),
( '-', 1),
( ' ', 1),
( ' ', 1),
( '💪', 1),
( '\U0001f90d', 1),
( '👉', 1),
( '😬', 1),
( '📺', 1),
( '🐮', 1),
( '🐷', 1),
( '*', 1),
( '🔪', 1),
( '💣', 1),
( '🐱', 1),
( '📺', 1),
( '🌿', 1),
( '🐸', 1),
( '—', 1),
( '🔥', 1),
( '🔒', 1)]
```

Top 3 special characters

```
1 res[:3] #including space
```

```
[(' ', 1801), ('#', 68), ('.', 67)]
```

```
1 res[1:4] #exluding space space
```

```
[('#', 68), ('.', 67), (',', 49)]
```