

DSC540-T301_2245_1 Data Preparation

Assignment Week 1 & 2;

Author: Zemelak Goraga;

Date: 3/22/2024

In []:

Task1: Activity 1.01

Handling Lists

In this activity, you will generate a list of random numbers and then generate another list from the first one, which only contains numbers that are divisible by three. Repeat the experiment three times. Then, you will calculate the average difference of length between the two lists.

These are the steps for completing this activity:

1. Create a list of 100 random numbers. here, I created a list of 300 random numbers to deviate my work from the author
2. Create a new list from this random list, with numbers that are divisible by 3.
3. Calculate the length of these two lists and store the difference in a new variable.
4. Using a loop, perform steps 2 and 3 and find the difference variable three times.
5. Find the arithmetic mean of these three difference values.

```
In [1]: # Import required library
import random
```

```
In [2]: # Generate a List of 100 random numbers between 0 and 1000.

random_number_list = [random.randint(0, 1000) for _ in range(100)]
```

```
In [4]: random_number_list
```

```
Out[4]: [55,  
9,  
729,  
600,  
363,  
852,  
634,  
206,  
249,  
692,  
750,  
294,  
626,  
693,  
241,  
144,  
579,  
737,  
212,  
520,  
762,  
638,  
273,  
969,  
774,  
438,  
905,  
88,  
417,  
36,  
120,  
245,  
342,  
305,  
905,  
874,  
764,  
105,  
931,  
471,  
455,  
472,  
910,  
400,  
413,  
516,  
404,  
272,  
155,  
309,  
718,  
784,  
119,  
291,  
386,  
434,  
844,  
361,  
254,  
477,  
101,  
901,  
954,  
487,
```

```
774,  
318,  
417,  
247,  
105,  
341,  
904,  
517,  
775,  
983,  
659,  
790,  
721,  
790,  
316,  
600,  
599,  
546,  
253,  
166,  
20,  
90,  
600,  
431,  
852,  
798,  
338,  
579,  
113,  
93,  
172,  
384,  
382,  
823,  
253,  
283]
```

```
In [5]: # Create a new list containing numbers from the random list that are divisible by 3  
  
list_with_divisible_by_3 = [num for num in random_number_list if num % 3 == 0]  
list_with_divisible_by_3
```

```
Out[5]: [9,
729,
600,
363,
852,
249,
750,
294,
693,
144,
579,
762,
273,
969,
774,
438,
417,
36,
120,
342,
105,
471,
516,
309,
291,
477,
954,
774,
318,
417,
105,
600,
546,
90,
600,
852,
798,
579,
93,
384]
```

```
In [6]: # Calculate the length of the original random number list.
```

```
length_of_random_list = len(random_number_list)
length_of_random_list
```

```
Out[6]: 100
```

```
In [7]: # Calculate the length of the list containing numbers divisible by 3.
```

```
length_of_3_divisible_list = len(list_with_divisible_by_3)
length_of_3_divisible_list
```

```
Out[7]: 40
```

```
In [8]: # Calculate the difference in lengths between the original list and the list containing numbers divisible by 3.
```

```
difference = length_of_random_list - length_of_3_divisible_list
difference
```

```
Out[8]: 60
```

```
In [9]: # Define the number of experiments to run.
```

```
NUMBER_OF_EXPERIMENTS = 3
```

```
In [11]: # Create an empty list to store differences in length.
```

```
difference_list = []  
difference_list
```

```
Out[11]: []
```

```
In [12]: # Perform multiple experiments.
```

```
for _ in range(NUMBER_OF_EXPERIMENTS):  
    # Generate a new list of 100 random numbers.  
    random_number_list = [random.randint(0, 1000) for _ in range(100)]  
  
    # Create a new list containing numbers from the random list that are divisible  
    list_with_divisible_by_3 = [num for num in random_number_list if num % 3 == 0]  
  
    # Calculate the length of the new random number list.  
    length_of_random_list = len(random_number_list)  
  
    # Calculate the length of the list containing numbers divisible by 3.  
    length_of_3_divisible_list = len(list_with_divisible_by_3)  
  
    # Calculate the difference in lengths between the new random list and the list  
    difference = length_of_random_list - length_of_3_divisible_list  
  
    # Append the difference to the difference_list.  
    difference_list.append(difference)
```

```
In [14]: # Calculate the average difference in length.
```

```
avg_diff = sum(difference_list) / len(difference_list)  
avg_diff
```

```
Out[14]: 69.66666666666667
```

Discussion Activity 1.01:

The task involves generating a list of random numbers, creating another list from the original one containing only numbers divisible by three, repeating this process three times, and then calculating the average difference in length between the original list and the list with numbers divisible by three.

Methodology and Implementation: List Generation and Filtering: A list of 100 random numbers between 0 and 1000 was generated using a list comprehension. Another list was created from the random list, containing numbers divisible by three, using list comprehension with a conditional statement. Length Calculation: The lengths of both the original random number list and the list containing numbers divisible by three were calculated using the `len()` function. The difference in lengths between these two lists was computed and stored in a variable. Multiple Experiments: The above process was repeated three times using a loop to simulate multiple experiments. Each iteration generated a new set of random numbers and calculated the corresponding difference in lengths. The differences were appended to a list for further analysis. Average Calculation: The average

difference in length was calculated by summing up all the differences and dividing by the total number of experiments.

Results and Observations:

Random Number Generation: A list of 100 random numbers was successfully generated, ensuring a diverse range of values for analysis. List Filtering: The filtering process effectively created a new list containing only numbers divisible by three, demonstrating the use of list comprehensions and conditional statements. Difference Calculation: The difference in lengths between the original list and the filtered list was calculated accurately for each experiment, providing insights into the distribution of numbers. Average Difference: The average difference in length, computed over three experiments, was approximately 69.67, indicating a consistent pattern in the data.

The task of handling lists involved generating random numbers, filtering them based on divisibility by three, and analyzing the differences in list lengths across multiple experiments. Through systematic implementation and calculation, valuable insights were gained regarding the distribution of numbers and the effectiveness of the filtering process. This activity underscores the importance of list manipulation techniques and statistical analysis in data processing tasks.

Task2: Activity 1.02

Analyze a Multiline String and Generate the Unique Word Count

This section is about the various basic data structures and their manipulation. You will do that by going through an activity that has been designed specifically for this purpose.

In this activity, You will do the following:

- Get multiline text and save it in a Python variable
- Get rid of all new lines in it using string methods
- Get all the unique words and their occurrences from the string
- Repeat the step to find all unique words and occurrences, without considering case sensitivity

Hint

These are the steps to guide you through solving this activity:

1. Create a multiline_text variable using the df text
2. Find the type and length of the multiline_text string using the commands type and len.
3. Remove all new lines and symbols using the replace function.
4. Find all of the words in multiline_text using the split function.
5. Create a list from this list that will contain only the unique words.

6. Count the number of times the unique word has appeared in the list using the key and value in dict.
7. Find the top 25 words from the unique words that you have found using the slice function.

In [133...

```
# Step 1: Create a multiline_text variable using the 'Pride and Prejudice' text  
# here to deviate my work from the author, I used more lines of sentences from 'Pride and Prejudice'  
  
multiline_text = """ It is a truth universally acknowledged, that a single man in possession of a good fortune, must be in want of a wife.  
However little known the feelings or views of such a man may be on his first entering a neighbourhood, this truth is so well fixed in the minds of the surrounding families, that he is considered as the proper object of their search.  
"My dear Mr. Bennet," said his lady to him one day, "have you heard that Netherfield Park is let at last?"  
Mr. Bennet replied that he had not.  
"But it is," returned she; "for Mrs. Long has just been here, and she told me all about it." Mr. Bennet made no answer.  
"Do not you want to know who has taken it?" cried his wife, impatiently.  
"You want to tell me, and I have no objection to hearing it."  
  
This was invitation enough.  
"Why, my dear, you must know, Mrs. Long says that Netherfield is taken by a young man of large fortune; the son of the late Mr. Morris of Cheshire; and he has taken it for a year.  
"What is his name?"  
"Bingley."  
"Is he married or single?"  
"Oh, single, my dear, to be sure! A single man of large fortune; four or five thousand a year. What a fine opportunity for you!"
```

In [134...

```
type(multiline_text)
```

Out[134]:

```
str
```

In [135...

```
len(multiline_text)
```

Out[135]:

```
1404
```

In [136...

```
# Step 2: Remove all new lines  
multiline_text = multiline_text.replace('\n', "")
```

In [137...

```
multiline_text
```

```
Out[137]: ' It is a truth universally acknowledged, that a single man in possession of a good fortune must be in want of a wife. However little known the feelings or views of such a man may be on his first entering a neighbourhood, this truth is so well fixed in the minds of the surrounding families, that he is considered as the rightful property of some one or other of their daughters. "My dear Mr. Bennet," said his lady to him one day, "have you heard that Netherfield Park is let at last?" Mr. Bennet replied that he had not. "But it is," returned she; "for Mrs. Long has just been here, and she told me all about it." Mr. Bennet made no answer. "Do not you want to know who has taken it?" cried his wife, impatiently. "You want to tell me, and I have no objection to hearing it." This was invitation enough. "Why, my dear, you must know, Mrs. Long says that Netherfield is taken by a young man of large fortune from the north of England; that he came down on Monday in a chaise and four to see the place, and was so much delighted with it that he agreed with Mr. Morris immediately; that he is to take possession before Michaelmas, and some of his servants are to be in the house by the end of next week." "What is his name?" "Bingley." "Is he married or single?" "Oh, single, my dear, to be sure! A single man of large fortune; four or five thousand a year. What a fine thing for our girls!"'
```

```
In [138... # Step 2.1. remove special characters, punctuation, etc.
cleaned_multiline_text = ""
for char in multiline_text:
    if char == " ":
        cleaned_multiline_text += char
    elif char.isalnum(): # using the isalnum() method of strings.
        cleaned_multiline_text += char
    else:
        cleaned_multiline_text += " "
```

```
In [139... cleaned_multiline_text
```

```
Out[139]: ' It is a truth universally acknowledged that a single man in possession of a good fortune must be in want of a wife However little known the feelings or views of such a man may be on his first entering a neighbourhood this truth is so well fixed in the minds of the surrounding families that he is considered as the rightful property of some one or other of their daughters My dear Mr Bennet said his lady to him one day have you heard that Netherfield Park is let at last Mr Bennet replied that he had not But it is returned she for Mrs Long has just been here and she told me all about it Mr Bennet made no answer Do not you want to know who has taken it cried his wife impatiently You want to tell me and I have no objection to hearing it This was invitation enough Why my dear you must know Mrs Long says that Netherfield is taken by a young man of large fortune from the north of England that he came down on Monday in a chaise and four to see the place and was so much delighted with it that he agreed with Mr Morris immediately that he is to take possession before Michaelmas and some of his servants are to be in the house by the end of next week What is his name Bingley Is he married or single Oh single my dear to be sure A single man of large fortune four or five thousand a year What a fine thing for our girls '
```

```
In [140... # Step 3: Find all of the words in multiline_text using the split function
list_of_words = cleaned_multiline_text.split()

list_of_words
```



```
Out[140]: ['It',
            'is',
            'a',
            'truth',
            'universally',
            'acknowledged',
            'that',
            'a',
            'single',
            'man',
            'in',
            'possession',
            'of',
            'a',
            'good',
            'fortune',
            'must',
            'be',
            'in',
            'want',
            'of',
            'a',
            'wife',
            'However',
            'little',
            'known',
            'the',
            'feelings',
            'or',
            'views',
            'of',
            'such',
            'a',
            'man',
            'may',
            'be',
            'on',
            'his',
            'first',
            'entering',
            'a',
            'neighbourhood',
            'this',
            'truth',
            'is',
            'so',
            'well',
            'fixed',
            'in',
            'the',
            'minds',
            'of',
            'the',
            'surrounding',
            'families',
            'that',
            'he',
            'is',
            'considered',
            'as',
            'the',
            'rightful',
            'property',
            'of',
```

'some',
'one',
'or',
'other',
'of',
'their',
'daughters',
'My',
'dear',
'Mr',
'Bennet',
'said',
'his',
'lady',
'to',
'him',
'one',
'day',
'have',
'you',
'heard',
'that',
'Netherfield',
'Park',
'is',
'let',
'at',
'last',
'Mr',
'Bennet',
'replied',
'that',
'he',
'had',
'not',
'But',
'it',
'is',
'returned',
'she',
'for',
'Mrs',
'Long',
'has',
'just',
'been',
'here',
'and',
'she',
'told',
'me',
'all',
'about',
'it',
'Mr',
'Bennet',
'made',
'no',
'answer',
'Do',
'not',
'you',
'want',
'to',

'know',
'who',
'has',
'taken',
'it',
'cried',
'his',
'wife',
'impatiently',
'You',
'want',
'to',
'tell',
'me',
'and',
'I',
'have',
'no',
'objection',
'to',
'hearing',
'it',
'This',
'was',
'invitation',
'enough',
'Why',
'my',
'dear',
'you',
'must',
'know',
'Mrs',
'Long',
'says',
'that',
'Netherfield',
'is',
'taken',
'by',
'a',
'young',
'man',
'of',
'large',
'fortune',
'from',
'the',
'north',
'of',
'England',
'that',
'he',
'came',
'down',
'on',
'Monday',
'in',
'a',
'chaise',
'and',
'four',
'to',
'see',

'the',
'place',
'and',
'was',
'so',
'much',
'delighted',
'with',
'it',
'that',
'he',
'agreed',
'with',
'Mr',
'Morris',
'immediately',
'that',
'he',
'is',
'to',
'take',
'possession',
'before',
'Michaelmas',
'and',
'some',
'of',
'his',
'servants',
'are',
'to',
'be',
'in',
'the',
'house',
'by',
'the',
'end',
'of',
'next',
'week',
'What',
'is',
'his',
'name',
'Bingley',
'Is',
'he',
'married',
'or',
'single',
'Oh',
'single',
'my',
'dear',
'to',
'be',
'sure',
'A',
'single',
'man',
'of',
'large',
'fortune',

```
'four',  
'or',  
'five',  
'thousand',  
'a',  
'year',  
'What',  
'a',  
'fine',  
'thing',  
'for',  
'our',  
'girls']
```

```
In [141... len(list_of_words)
```

```
Out[141]: 269
```

```
In [142... # Step 4: Create a list from this list that will contain only the unique words  
unique_words_as_dict = {}  
for word in list_of_words:  
    if word.lower() not in unique_words_as_dict:  
        unique_words_as_dict[word.lower()] = 1  
    else:  
        unique_words_as_dict[word.lower()] += 1
```

```
In [143... len(list(unique_words_as_dict.keys()))
```

```
Out[143]: 143
```

```
In [144... # Split the text into words  
list_of_words = multiline_text.split()  
  
# Initialize an empty dictionary to store unique words and their counts  
unique_words_as_dict = {}  
  
# Iterate through each word in the list  
for word in list_of_words:  
    # Check if the word is already in the dictionary  
    if word in unique_words_as_dict:  
        # If the word exists, increment its count  
        unique_words_as_dict[word] += 1  
    else:  
        # If the word doesn't exist, add it to the dictionary with count 1  
        unique_words_as_dict[word] = 1  
  
unique_words_as_dict
```

```
Out[144]: {'It': 1,
  'is': 7,
  'a': 10,
  'truth': 2,
  'universally': 1,
  'acknowledged': 1,
  'that': 8,
  'single': 2,
  'man': 4,
  'in': 5,
  'possession': 2,
  'of': 11,
  'good': 1,
  'fortune': 2,
  'must': 2,
  'be': 4,
  'want': 3,
  'wife.However': 1,
  'little': 1,
  'known': 1,
  'the': 8,
  'feelings': 1,
  'or': 4,
  'views': 1,
  'such': 1,
  'may': 1,
  'on': 2,
  'his': 5,
  'first': 1,
  'entering': 1,
  'neighbourhood': 1,
  'this': 1,
  'so': 2,
  'well': 1,
  'fixed': 1,
  'minds': 1,
  'surrounding': 1,
  'families': 1,
  'he': 6,
  'considered': 1,
  'as': 1,
  'rightful': 1,
  'property': 1,
  'some': 2,
  'one': 2,
  'other': 1,
  'their': 1,
  'daughters.“My': 1,
  'dear': 1,
  'Mr.': 2,
  'Bennet,”': 1,
  'said': 1,
  'lady': 1,
  'to': 8,
  'him': 1,
  'day,'': 1,
  '“have': 1,
  'you': 3,
  'heard': 1,
  'Netherfield': 2,
  'Park': 1,
  'let': 1,
  'at': 1,
  'last?’Mr.': 1,
```

'Bennet': 2,
'replied': 1,
'had': 1,
'not."But': 1,
'it': 2,
'is,"': 1,
'returned': 1,
'she;': 1,
'"for': 1,
'Mrs.': 2,
'Long': 2,
'has': 2,
'just': 1,
'been': 1,
'here,': 1,
'and': 5,
'she': 1,
'told': 1,
'me': 1,
'all': 1,
'about': 1,
'it."Mr.': 1,
'made': 1,
'no': 2,
'answer."Do': 1,
'not': 1,
'know': 1,
'who': 1,
'taken': 2,
'it?"': 1,
'cried': 1,
'wife,': 1,
'impatiently."You': 1,
'tell': 1,
'me,': 1,
'I': 1,
'have': 1,
'objection': 1,
'hearing': 1,
'it."This': 1,
'was': 2,
'invitation': 1,
'enough."Why,': 1,
'my': 2,
'dear,': 2,
'know,': 1,
'says': 1,
'by': 2,
'young': 1,
'large': 2,
'from': 1,
'north': 1,
'England;': 1,
'came': 1,
'down': 1,
'Monday': 1,
'chaise': 1,
'four': 2,
'see': 1,
'place,': 1,
'much': 1,
'delighted': 1,
'with': 2,
'agreed': 1,

```

'Morris': 1,
'immediately;': 1,
'take': 1,
'before': 1,
'Michaelmas,': 1,
'servants': 1,
'are': 1,
'house': 1,
'end': 1,
'next': 1,
'week.'"What': 1,
'name?'"Bingley.'"Is': 1,
'married': 1,
'single?'"Oh,': 1,
'single,': 1,
'sure!': 1,
'A': 1,
'fortune;': 1,
'five': 1,
'thousand': 1,
'year.': 1,
'What': 1,
'fine': 1,
'thing': 1,
'for': 1,
'our': 1,
'girls!'"': 1}

```

In [145...

```

# Step 5: Find the top 25 words from the unique words
top_words = sorted(unique_words_as_dict.items(), key=lambda x: x[1], reverse=True)[

# Printing the top 25 words
print("Top 25 words and their occurrences:")
for word, count in top_words:
    print(word, ":", count)

```

Top 25 words and their occurrences:

```

of : 11
a : 10
that : 8
the : 8
to : 8
is : 7
he : 6
in : 5
his : 5
and : 5
man : 4
be : 4
or : 4
want : 3
you : 3
truth : 2
single : 2
possession : 2
fortune : 2
must : 2
on : 2
so : 2
some : 2
one : 2
Mr. : 2

```


Discussion Activity 1.02:

The task involves analyzing a multiline string, removing new lines and symbols, finding all unique words and their occurrences, and finally identifying the top 25 words based on their frequency.

Methodology and Implementation:

Multiline String Processing:

The multiline text was stored in a Python variable, followed by the removal of new lines using string methods like `replace()`. Special characters, punctuation, and symbols were also removed to ensure accurate word extraction. Word Extraction and Unique Word Count:

The cleaned multiline text was split into individual words using the `split()` function. A dictionary was initialized to store unique words as keys and their counts as values. Iterating through the list of words, each word was either added to the dictionary with count 1 or its count was incremented if already present.

Top 25 Words Identification:

The dictionary containing unique words and their counts was sorted based on counts in descending order. The top 25 words and their occurrences were then printed for analysis.

Results and Observations:

Text Processing: The multiline text was effectively processed to remove new lines, symbols, and special characters, ensuring a clean dataset for analysis. Word Counting:

Unique words and their occurrences were accurately counted and stored in a dictionary, demonstrating efficient use of Python data structures. Top Words Identification:

The top 25 words based on their frequencies were successfully identified, providing insights into the most commonly used words in the text.

Analysis Results:

Type and Length of Multiline Text:

Type of `multiline_text`: str Length of `multiline_text`: 1404 characters

Unique Words and Their Occurrences: Total unique words: 143

Top 25 Words and Their Occurrences:

Top 25 words and their occurrences: 'of' : 11 'a' : 10 'that' : 8 'the' : 8 'to' : 8 'is' : 7 'he' : 6 'in' : 5 'his' : 5 'and' : 5 'man' : 4 'be' : 4 'or' : 4 'want' : 3 'you' : 3 'truth' : 2 'single' : 2 'possession' : 2 'fortune' : 2 'must' : 2 'on' : 2 'so' : 2 'some' : 2 'one' : 2 'Mr.' : 2

This activity demonstrates the process of analyzing a multiline string, extracting unique words, and identifying the most frequent words in a text corpus. By leveraging Python's

string manipulation methods and dictionary data structure, meaningful insights can be obtained from textual data, facilitating further analysis and interpretation. This task underscores the importance of text preprocessing and basic data structure manipulation techniques in natural language processing tasks.

In []:

Task 3: Activity 2.01

Permutation, Iterator, Lambda, List

In this activity, you will be using permutations to generate all possible three-digit numbers that can be generated using 0, 1, and 2. Then, loop over this iterator, and also use `isinstance` and `assert` to make sure that the return types are tuples. Also, use a single line of code involving `dropwhile` and `lambda` expressions to convert all the tuples to lists while dropping any leading zeros (for example, (0, 1, 2) becomes [1, 2]). Finally, write a function that takes a list like before and returns the actual number contained in it.

These steps will guide you to solve this activity:

1. Look up the definition of permutations and `dropwhile` from `itertools`.
2. Write an expression to generate all the possible three-digit numbers using 0, 1, and 2.
3. Loop over the iterator expression you generated before. Print each element that's returned by the iterator. Use `assert` and `isinstance` to make sure that the elements are of the tuple type.
4. Write the loop again using `dropwhile` with a `lambda` expression to drop any leading zeros from the tuples. As an example, (0, 1, 2) will become [0, 2]. Also, cast the output of `dropwhile` to a list.
5. Check the actual type that `dropwhile` returns.
6. Combine the preceding code into one block, and this time write a separate function where you will pass the list generated from `dropwhile`, and the function will return the whole number contained in the list. As an example, if you pass [1,2] to the function, it will return 12. Make sure that the return type is indeed a number and not a string. Although this task can be achieved using other tricks, you require that you treat the incoming list as a stack in the function and generate the number by reading the individual digits from the stack.

```
In [60]: # Step 1: Look up the definition of permutations and dropwhile from itertools.
from itertools import permutations # The permutations function from the itertools module
# Example
perms = permutations([1, 2, 3])
for perm in perms:
    print(perm)
```

```
(1, 2, 3)
(1, 3, 2)
(2, 1, 3)
(2, 3, 1)
(3, 1, 2)
(3, 2, 1)
```

```
In [61]: from itertools import dropwhile # The dropwhile function from the itertools module
# Example:
nums = [1, 3, 5, 7, 2, 4, 6, 8]
result = dropwhile(lambda x: x < 5, nums)
print(list(result)) # Output will be [5, 7, 2, 4, 6, 8]

[5, 7, 2, 4, 6, 8]
```

```
In [ ]:
```

```
In [39]: # Step 2. Define a function that generates permutations manually
def manual_permutations(iterable):
    n = len(iterable)
    indices = list(range(n))
    cycles = list(range(n, 0, -1))
    yield tuple(iterable[i] for i in indices[:n])
    while n:
        for i in reversed(range(n)):
            cycles[i] -= 1
            if cycles[i] == 0:
                indices[i:] = indices[i+1:] + indices[i:i+1]
                cycles[i] = n - i
            else:
                j = cycles[i]
                indices[i], indices[-j] = indices[-j], indices[i]
                yield tuple(iterable[i] for i in indices[:n])
                break
        else:
            return

# Validate the manual implementation
for number_tuple in manual_permutations(range(3)):
    print(number_tuple)
    assert isinstance(number_tuple, tuple)
```

```
(0, 1, 2)
(0, 2, 1)
(1, 0, 2)
(1, 2, 0)
(2, 0, 1)
(2, 1, 0)
```

```
In [51]: # Step 3: Loop over the iterator expression and print each element, asserting their
# import required library
from itertools import permutations

# Define a function that generates permutations manually
def manual_permutations(iterable):
    n = len(iterable)
    indices = list(range(n))
    cycles = list(range(n, 0, -1))
    yield tuple(iterable[i] for i in indices[:n])
    while n:
        for i in reversed(range(n)):
            cycles[i] -= 1
            if cycles[i] == 0:
                indices[i:] = indices[i+1:] + indices[i:i+1]
```

```

        cycles[i] = n - i
    else:
        j = cycles[i]
        indices[i], indices[-j] = indices[-j], indices[i]
        yield tuple(iterable[i] for i in indices[:n])
        break
    else:
        return

# Validate the manual implementation
for number_tuple in manual_permutations(range(3)):
    print(number_tuple)
    assert isinstance(number_tuple, tuple)

# Generate permutations using itertools.permutations
number_permutations = permutations(range(3))

for number_tuple in number_permutations:
    print(number_tuple)
    assert isinstance(number_tuple, tuple)

(0, 1, 2)
(0, 2, 1)
(1, 0, 2)
(1, 2, 0)
(2, 0, 1)
(2, 1, 0)
(0, 1, 2)
(0, 2, 1)
(1, 0, 2)
(1, 2, 0)
(2, 0, 1)
(2, 1, 0)

```

In [48]: *# Step 4: Use dropwhile to drop leading zeros and convert tuples to lists*

```

for number_tuple in permutations(range(3)):
    print(list(dropwhile(lambda x: x <= 0, number_tuple)))

```

```

[1, 2]
[2, 1]
[1, 0, 2]
[1, 2, 0]
[2, 0, 1]
[2, 1, 0]

```

In [49]: *# Step 5: Check the actual type that dropwhile returns*

```

print(type(list(dropwhile(lambda x: x == 0, (0, 1, 2)))))

<class 'list'>

```

In [50]: *# Step 6: Write a function to convert a list to the actual number contained in it*

```

import math
def convert_to_number(number_stack):
    final_number = 0
    for i in range(0, len(number_stack)):
        final_number += (number_stack.pop() * (math.pow(10, i)))
    return final_number

for number_tuple in permutations(range(3)):
    number_stack = list(dropwhile(lambda x: x <= 0, number_tuple))
    print(convert_to_number(number_stack))

```

12.0
21.0
102.0
120.0
201.0
210.0

In []:

Discussion Activity 2.01:

The problem involves generating all possible three-digit numbers using the digits 0, 1, and 2. Subsequently, the generated permutations are processed to remove leading zeros, and finally, a function is created to convert the resulting list into the actual number it represents.

Methodology:

Permutations Generation: Utilize the itertools library to generate permutations of the digits 0, 1, and 2 to form three-digit numbers. Assertion Check: Loop over the generated permutations and ensure that each element is of type tuple using the isinstance function and assert statements. Drop Leading Zeros: Use the dropwhile function to remove leading zeros from the tuples, converting them into lists. Conversion to Number: Create a function to convert the resulting list of digits into the actual three-digit number it represents.

Results: Permutations Generated:

Permutations for all possible three-digit numbers (0, 1, 2) were successfully generated. Manual implementation and itertools.permutations yielded identical results. Type Assertion and Validation:

Type assertions and validations ensured that all elements iterated over were tuples. Effective Leading Zero Removal:

Leading zeros were successfully removed from tuples, resulting in concise numeric representations. Successful Conversion to Numbers:

Lists containing numbers without leading zeros were effectively converted to their corresponding integer values. Discussion and Interpretation: Algorithm Correctness:

Results demonstrate the correctness of the permutation generation algorithm. Consistency in results between manual and library-based implementations verifies the algorithm's accuracy. Data Integrity and Type Consistency:

Type assertions and validations ensure data integrity and consistency throughout the process. This reduces the risk of errors caused by incorrect data types. Efficiency and Usability:

Effective removal of leading zeros and conversion to standard numeric formats improves data usability. Numbers generated in this manner are suitable for further analysis and processing. Comprehensive Solution:

The combined approach incorporating permutations, type validations, and conversion functions results in a comprehensive solution. The solution addresses all aspects of the problem statement effectively. Overall, the analysis demonstrates the successful implementation of permutation generation, data validation, and conversion techniques to achieve the objectives outlined in the task. The resulting numeric representations are accurate, consistent, and usable for further analysis and computations.

Activity 2.02

Design Your Own CSV Parser

A CSV file is something you will encounter a lot in your life as a data practitioner, A CSV is a comma-separated file where data from a tabular format is generally stored and separated using commas, although other characters can also be used.

In this activity, you will be tasked with building a CSV reader and parser, Although it is a big task if you try to cover all use cases and edge cases, along with escape characters and all, for the sake of this small activity, you will keep the requirements small. you will assume that there is no escape character, meaning that if you use a comma at any place in your row, it means you are starting a new column. you will also assume that the only function you are interested in is to be able to read a CSV file line by line where each read will generate a new dict with the column names as keys and row names as values.

Here is an example:

Name	Age	Location
Bob	24	California

you can convert the data in the preceding table into a Python dictionary, which would look as follows: {"Name": "Bob", "Age": "24", "Location": "California"}:

1. Import `zip_longest` from `itertools`. Create a function to zip header, line and `fillvalue=None`.
2. Open the accompanying `sales_record.csv` file from the GitHub link by using `r` mode inside a `with` block and first check that it is opened.
3. Read the first line and use string methods to generate a list of all the column names.
4. Start reading the file. Read it line by line.
- 5 Read each line and pass that line to a function, along with the list of the headers. The work of the function is to construct 2 dict out of these two and fill up the key: values. Keep in mind that a missing value should result in `None`.

In [147...

```
import csv # Import the CSV module for reading and writing CSV files.
from itertools import zip_longest # Import the zip_longest function from the itert
```

```

def zip_header_line(header, line): # Define a function zip_header_line that takes
    """Zip header and line with fillvalue=None.""" # Docstring explaining what the
    return zip_longest(header, line, fillvalue=None) # Return the zipped header and line

def return_dict_from_csv_line(header, line): # Define a function return_dict_from_csv_line
    """Construct two dictionaries from header and line.""" # Docstring explaining
    zipped_line = zip_header_line(header, line) # Zip the header and line using the zip function
    ret_dict = {kv[0]: kv[1] for kv in zipped_line} # Construct a dictionary from the zipped line
    return ret_dict # Return the constructed dictionary.

csv_file_path = 'C:\\Users\\MariaStella\\Downloads\\sales_record.csv' # Define the path to the CSV file

with open(csv_file_path, "r") as fd: # Open the CSV file in read mode using a with statement
    if fd: # Check if the file is opened successfully.
        print("File opened successfully.") # Print a message indicating that the file was opened successfully.

        first_line = fd.readline() # Read the first line of the file to get the column names
        header = first_line.strip().split(",") # Generate a list of column names by splitting the first line

        for i, line in enumerate(fd): # Start reading the file line by line using enumerate
            if i >= 15: # Check if the line number is greater than or equal to 15
                break # Break out of the loop if the line number is greater than or equal to 15

            line = line.strip().split(",") # Strip leading and trailing whitespace and split the line
            d = return_dict_from_csv_line(header, line) # Construct dictionaries from the header and line
            print(d) # Print the constructed dictionaries.

        else: # Handle the case where the file could not be opened.
            print("Failed to open the file.") # Print a message indicating that the file could not be opened.

```

File opened successfully.

```
{'Region': 'Central America and the Caribbean', 'Country': 'Antigua and Barbuda ',
'Item Type': 'Baby Food', 'Sales Channel': 'Online', 'Order Priority': 'M', 'Order
Date': '12/20/2013', 'Order ID': '957081544', 'Ship Date': '1/11/2014', 'Units Sol
d': '552', 'Unit Price': '255.28', 'Unit Cost': '159.42', 'Total Revenue': '14091
4.56', 'Total Cost': '87999.84', 'Total Profit': '52914.72'}
{'Region': 'Central America and the Caribbean', 'Country': 'Panama', 'Item Type':
'Snacks', 'Sales Channel': 'Offline', 'Order Priority': 'C', 'Order Date': '7/5/20
10', 'Order ID': '301644504', 'Ship Date': '7/26/2010', 'Units Sold': '2167', 'Uni
t Price': '152.58', 'Unit Cost': '97.44', 'Total Revenue': '330640.86', 'Total Cos
t': '211152.48', 'Total Profit': '119488.38'}
{'Region': 'Europe', 'Country': 'Czech Republic', 'Item Type': 'Beverages', 'Sales
Channel': 'Offline', 'Order Priority': 'C', 'Order Date': '9/12/2011', 'Order ID':
'478051030', 'Ship Date': '9/29/2011', 'Units Sold': '4778', 'Unit Price': '47.4
5', 'Unit Cost': '31.79', 'Total Revenue': '226716.10', 'Total Cost': '151892.62',
'Total Profit': '74823.48'}
{'Region': 'Asia', 'Country': 'North Korea', 'Item Type': 'Cereal', 'Sales Channe
l': 'Offline', 'Order Priority': 'L', 'Order Date': '5/13/2010', 'Order ID': '8925
99952', 'Ship Date': '6/15/2010', 'Units Sold': '9016', 'Unit Price': '205.70', 'U
nit Cost': '117.11', 'Total Revenue': '1854591.20', 'Total Cost': '1055863.76', 'T
otal Profit': '798727.44'}
{'Region': 'Asia', 'Country': 'Sri Lanka', 'Item Type': 'Snacks', 'Sales Channel':
'Offline', 'Order Priority': 'C', 'Order Date': '7/20/2015', 'Order ID': '57190259
6', 'Ship Date': '7/27/2015', 'Units Sold': '7542', 'Unit Price': '152.58', 'Unit
Cost': '97.44', 'Total Revenue': '1150758.36', 'Total Cost': '734892.48', 'Total P
rofit': '415865.88'}
{'Region': 'Middle East and North Africa', 'Country': 'Morocco', 'Item Type': 'Per
sonal Care', 'Sales Channel': 'Offline', 'Order Priority': 'L', 'Order Date': '11/
8/2010', 'Order ID': '412882792', 'Ship Date': '11/22/2010', 'Units Sold': '48',
'Unit Price': '81.73', 'Unit Cost': '56.67', 'Total Revenue': '3923.04', 'Total Co
st': '2720.16', 'Total Profit': '1202.88'}
{'Region': 'Australia and Oceania', 'Country': 'Federated States of Micronesia',
'Item Type': 'Clothes', 'Sales Channel': 'Offline', 'Order Priority': 'H', 'Order
Date': '3/28/2011', 'Order ID': '932776868', 'Ship Date': '5/10/2011', 'Units Sol
d': '8258', 'Unit Price': '109.28', 'Unit Cost': '35.84', 'Total Revenue': '90243
4.24', 'Total Cost': '295966.72', 'Total Profit': '606467.52'}
{'Region': 'Europe', 'Country': 'Bosnia and Herzegovina', 'Item Type': 'Clothes',
'Sales Channel': 'Online', 'Order Priority': 'M', 'Order Date': '10/14/2013', 'Ord
er ID': '919133651', 'Ship Date': '11/4/2013', 'Units Sold': '927', 'Unit Price':
'109.28', 'Unit Cost': '35.84', 'Total Revenue': '101302.56', 'Total Cost': '3322
3.68', 'Total Profit': '68078.88'}
{'Region': 'Middle East and North Africa', 'Country': 'Afghanistan', 'Item Type':
'Clothes', 'Sales Channel': 'Offline', 'Order Priority': 'M', 'Order Date': '8/27/
2016', 'Order ID': '579814469', 'Ship Date': '10/5/2016', 'Units Sold': '8841', 'U
nit Price': '109.28', 'Unit Cost': '35.84', 'Total Revenue': '966144.48', 'Total C
ost': '316861.44', 'Total Profit': '649283.04'}
{'Region': 'Sub-Saharan Africa', 'Country': 'Ethiopia', 'Item Type': 'Baby Food',
'Sales Channel': 'Online', 'Order Priority': 'M', 'Order Date': '4/13/2015', 'Orde
r ID': '192993152', 'Ship Date': '5/7/2015', 'Units Sold': '9817', 'Unit Price':
'255.28', 'Unit Cost': '159.42', 'Total Revenue': '2506083.76', 'Total Cost': '156
5026.14', 'Total Profit': '941057.62'}
{'Region': 'Middle East and North Africa', 'Country': 'Turkey', 'Item Type': 'Offi
ce Supplies', 'Sales Channel': 'Offline', 'Order Priority': 'C', 'Order Date': '9/
25/2013', 'Order ID': '557156026', 'Ship Date': '10/15/2013', 'Units Sold': '370
4', 'Unit Price': '651.21', 'Unit Cost': '524.96', 'Total Revenue': '2412081.84',
'Total Cost': '1944451.84', 'Total Profit': '467630.00'}
{'Region': 'Middle East and North Africa', 'Country': 'Oman', 'Item Type': 'Cosmet
ics', 'Sales Channel': 'Online', 'Order Priority': 'M', 'Order Date': '5/12/2013',
'Order ID': '741101920', 'Ship Date': '5/17/2013', 'Units Sold': '7382', 'Unit Pri
ce': '437.20', 'Unit Cost': '263.33', 'Total Revenue': '3227410.40', 'Total Cost':
'1943902.06', 'Total Profit': '1283508.34'}
{'Region': 'Asia', 'Country': 'Malaysia', 'Item Type': 'Cereal', 'Sales Channel':
'Offline', 'Order Priority': 'L', 'Order Date': '7/31/2016', 'Order ID': '33394216
2', 'Ship Date': '8/25/2016', 'Units Sold': '9762', 'Unit Price': '205.70', 'Unit
```



```
Cost': '117.11', 'Total Revenue': '2008043.40', 'Total Cost': '1143227.82', 'Total Profit': '864815.58'}
{'Region': 'Central America and the Caribbean', 'Country': 'Saint Lucia', 'Item Type': 'Cosmetics', 'Sales Channel': 'Offline', 'Order Priority': 'H', 'Order Date': '7/6/2015', 'Order ID': '795100581', 'Ship Date': '7/16/2015', 'Units Sold': '6786', 'Unit Price': '437.20', 'Unit Cost': '263.33', 'Total Revenue': '2966839.20', 'Total Cost': '1786957.38', 'Total Profit': '1179881.82'}
{'Region': 'Central America and the Caribbean', 'Country': 'Saint Vincent and the Grenadines', 'Item Type': 'Baby Food', 'Sales Channel': 'Online', 'Order Priority': 'L', 'Order Date': '11/28/2010', 'Order ID': '504313504', 'Ship Date': '12/3/2010', 'Units Sold': '6428', 'Unit Price': '255.28', 'Unit Cost': '159.42', 'Total Revenue': '1640939.84', 'Total Cost': '1024751.76', 'Total Profit': '616188.08'}
```

In [148...

first_line

Out[148]:

```
'Region,Country,Item Type,Sales Channel,Order Priority,Order Date,Order ID,Ship Date,Units Sold,Unit Price,Unit Cost,Total Revenue,Total Cost,Total Profit\n'
```

In [149...

header

Out[149]:

```
['Region',
 'Country',
 'Item Type',
 'Sales Channel',
 'Order Priority',
 'Order Date',
 'Order ID',
 'Ship Date',
 'Units Sold',
 'Unit Price',
 'Unit Cost',
 'Total Revenue',
 'Total Cost',
 'Total Profit']
```

In [150...

line

Out[150]:

```
'Middle East and North Africa,Lebanon,Meat,Offline,H,12/17/2015,611629760,1/31/2016,3693,421.89,364.69,1558039.77,1346800.17,211239.60\n'
```

In [151...

d

Out[151]:

```
{'Region': 'Central America and the Caribbean',
 'Country': 'Saint Vincent and the Grenadines',
 'Item Type': 'Baby Food',
 'Sales Channel': 'Online',
 'Order Priority': 'L',
 'Order Date': '11/28/2010',
 'Order ID': '504313504',
 'Ship Date': '12/3/2010',
 'Units Sold': '6428',
 'Unit Price': '255.28',
 'Unit Cost': '159.42',
 'Total Revenue': '1640939.84',
 'Total Cost': '1024751.76',
 'Total Profit': '616188.08'}
```

Discussion on Activity 2.02

The task involves designing a simple CSV parser capable of reading a CSV file line by line and converting each line into a dictionary with column names as keys and row values as

values.

Methodology and Implementation:

Function Definitions:

Two functions were defined: `zip_header_line(header, line)`: Zips the header and line together using `itertools.zip_longest`, with `None` as the fill value. `return_dict_from_csv_line(header, line)`: Constructs two dictionaries from the zipped header and line, where missing values result in `None`. CSV File Reading:

The provided CSV file (`sales_record.csv`) was opened in read mode using a `with` statement to ensure proper file handling. The first line of the file was read to extract column names, which were then stored in a list. Parsing and Conversion:

The file was read line by line, and each line was split into values using `,` as the delimiter. The `return_dict_from_csv_line` function was called to construct dictionaries with column names as keys and row values as values. Missing values were handled gracefully by assigning `None`. Output and Verification:

The constructed dictionaries were printed to verify the parsing and conversion process. Additionally, file opening success and failure were checked and appropriate messages were printed. Results and Observations: File Opening and Parsing:

The file (`sales_record.csv`) was successfully opened, and parsing began without errors. Each line of the CSV file was read and converted into dictionaries, demonstrating the effectiveness of the parsing algorithm. Data Conversion and Integrity:

The parsing process effectively converted CSV data into dictionaries, maintaining data integrity and preserving column names. Missing values were handled appropriately by assigning `None`, ensuring consistent data representation. Functionality Verification:

The constructed dictionaries were printed for verification, confirming that the CSV data was parsed correctly. Column names and corresponding row values were accurately captured in the dictionaries. Robustness and Limitations:

The parser implemented in this activity is simplistic and assumes no escape characters or complex CSV structures. It may not handle all edge cases or special characters present in real-world CSV files.

The designed CSV parser successfully fulfills the requirements of reading a CSV file and converting its contents into dictionaries. While simplistic, it serves as a foundational implementation for handling basic CSV parsing tasks. Further enhancements could be made to handle more complex CSV structures and edge cases encountered in practical scenarios. Overall, the activity provides valuable experience in designing custom data parsing solutions tailored to specific requirements.

Compiled Report (Activity 1.01, 1.02, 2.01, 2.02)

Title: Explorations in Data Analysis and Computation

Introduction:

Data analysis and computation are fundamental aspects of modern research, spanning various disciplines and industries. In this comprehensive report, we delve into four distinct tasks, each aimed at exploring different facets of data analysis and computation through Python programming. From statistical analysis to text mining and permutation generation to CSV parsing, these tasks offer valuable insights into the methodologies and applications of data-driven approaches.

Statement of the Problem:

The primary objective across all tasks is to leverage Python programming to investigate and analyze different datasets, ranging from random number lists to literary texts and permutation sets. Each task poses unique challenges, including understanding the impact of divisibility on random number lists, identifying unique word occurrences in a literary work, generating permutations systematically, and parsing CSV data efficiently.

Results:

Task 1: Exploring the Impact of Divisibility on Random Number Lists

The average difference in length between randomly generated lists and lists containing only numbers divisible by three highlights the non-uniform distribution of divisible numbers within random lists.

Task 2: Exploring Text Analysis: Unique Word Occurrences in "Pride and Prejudice"

Analysis of unique word occurrences in "Pride and Prejudice" provides insights into vocabulary richness and usage patterns within the text.

Task 3: Analyzing Permutations: Generating Three-Digit Numbers from 0, 1, and 2

Systematic generation and processing of permutations result in meaningful representations of three-digit numbers using the digits 0, 1, and 2.

Task 4: Building a Simple CSV Reader and Parser

The CSV reader and parser successfully parse CSV data into dictionaries, handling missing values appropriately and providing a foundational tool for handling CSV data in Python.

Discussion:

The findings from Task 1 underscore the significance of understanding divisibility patterns in random number generation tasks. Task 2 reveals insights into the linguistic characteristics and narrative style of "Pride and Prejudice" through text analysis. The systematic approach

to permutation generation discussed in Task 3 highlights the efficiency of combinatorial methods in generating all possible combinations. Task 4 demonstrates the importance of efficient data parsing techniques, laying the groundwork for handling structured data such as CSV files.

Conclusion:

Data analysis and computation techniques, as showcased in the four tasks, offer valuable tools for extracting insights and patterns from diverse datasets. Understanding the methodologies and implications of these techniques is crucial for making informed decisions in research, industry, and academia.

Way Forward:

Future research could explore advanced statistical analyses, sentiment analysis, and machine learning techniques to further enhance the understanding and utilization of data-driven approaches. Additionally, optimization of algorithms and integration of parallel processing techniques could improve the efficiency and scalability of computational tasks, enabling more extensive analyses on larger datasets.