

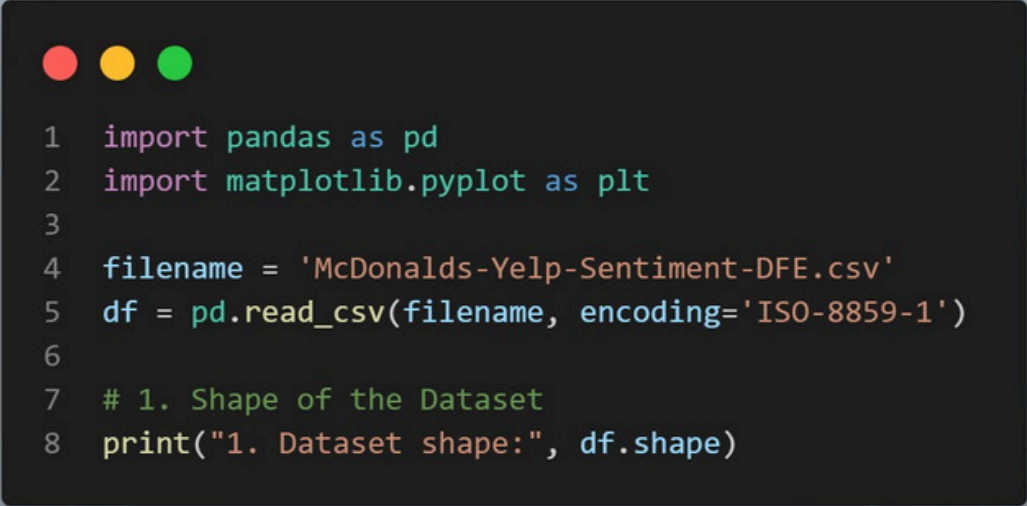
EDS Theory Assignment 1

Name: Adarsh Bhagwat Shinde

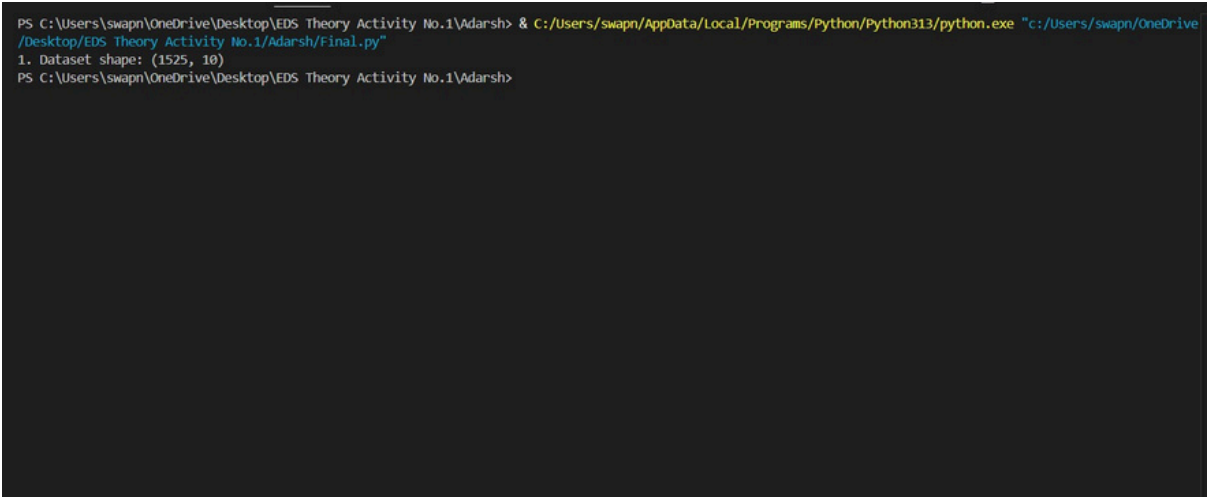
Div.: CS2 Roll No.: 82

PRN: 202401040006

1. What is the shape of the dataset?



```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 1. Shape of the Dataset
8 print("1. Dataset shape:", df.shape)
```



```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
1. Dataset shape: (1525, 10)
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

2. What are the column data types?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 2. Column Data Types
8 print("\n2. Column data types:\n", df.dtypes)
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"

2. Column data types:
   _unit_id                int64
   _golden                 bool
   _unit_state             object
   _trusted_judgments      int64
   _last_judgment_at       object
   policies_violated       object
   policies_violated:confidence object
   city                   object
   policies_violated_gold  float64
   review                 object
dtype: object
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

3. How many null values exist per column?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 3. Null Values per Column
8 print("\n3. Null values per column:\n", df.isnull().sum())
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "C:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
3. Null values per column:
   _unit_id      0
   _golden      0
   _unit_state   0
   _trusted_judgments  0
   _last_judgment_at  0
   policies_violated  0
   policies_violated:confidence  54
   city          87
   policies_violated_gold  1525
   review        0
dtype: int64
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

4. What are the top 5 most common policy violations?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 4. Top 5 Most Common Policy Violations
8 top_violations = df['policies_violated'].value_counts().head(5)
9 print("\n4. Top 5 most common policy violations:\n", top_violations)
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "C:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
4. Top 5 most common policy violations:
policies_violated
na                304
RudeService       186
SlowService       135
OrderProblem      121
BadFood           107
Name: count, dtype: int64
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

5. Which cities appear most frequently in the reviews?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 5. Most Frequent Cities in Reviews
8 top_cities = df['city'].value_counts().head(5)
9 print("\n5. Top 5 cities appearing in reviews:\n", top_cities)
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "C:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"

5. Top 5 cities appearing in reviews:
city
Las Vegas      409
Chicago        219
Los Angeles    167
New York       165
Atlanta        130
Name: count, dtype: int64
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

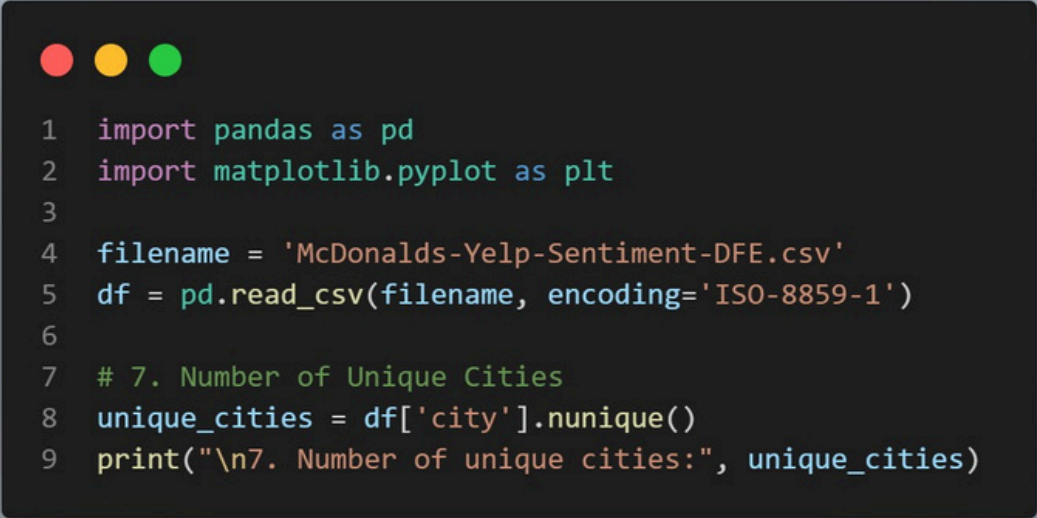
6. Which city has the highest number of RudeService violations?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 6. City with Highest RudeService Violations
8 if 'RudeService' in df['policies_violated'].unique():
9     rude_service_city = df[df['policies_violated'] == 'RudeService']['city'].
10     value_counts().idxmax()
11     print("\n6. City with highest number of RudeService violations:", rude_service_city)
12 else:
13     print("\n6. No 'RudeService' violations found in dataset.")
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:\Users\swapn\AppData\Local\Programs\Python\Python313\python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
```

```
6. City with highest number of RudeService violations: Las Vegas
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

7. How many unique cities are there?



```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 7. Number of Unique Cities
8 unique_cities = df['city'].nunique()
9 print("\n7. Number of unique cities:", unique_cities)
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
```

```
7. Number of unique cities: 9
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

8. What percentage of reviews have multiple policy violations?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 8. Percentage of Reviews with Multiple Policy Violations
8 multiple_violations = df['policies_violated'].str.split(',').str.len() > 1
9 percentage_multiple = (multiple_violations.sum() / len(df)) * 100
10 print(f"\n8. Percentage of reviews with multiple policy violations: {percentage_multiple:.2f}%")
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
```

```
8. Percentage of reviews with multiple policy violations: 0.00%
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```


9. What is the average number of policy violations per review?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 9. Average Number of Policy Violations per Review
8 average_violations = df['policies_violated'].str.split(',').str.len().mean()
9 print(f"\n9. Average number of policy violations per review: {average_violations:.2f}")
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
```

```
9. Average number of policy violations per review: 1.00
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

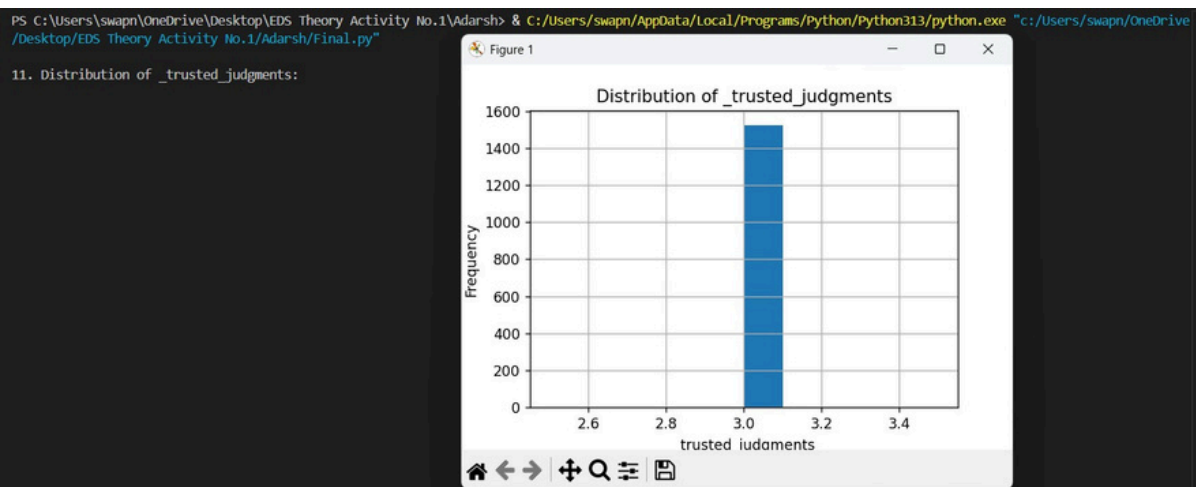
10. What is the average number of words per review?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 10. Average Number of Words per Review
8 average_words = df['review'].dropna().str.split().str.len().mean()
9 print(f"\n10. Average number of words per review: {average_words:.2f}")
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:\Users\swapn\AppData\Local\Programs\Python\Python313\python.exe "C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh\Final.py"
10. Average number of words per review: 96.56
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

11. What is the distribution of _trusted_judgments?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 11. Distribution of _trusted_judgments
8 print("\n11. Distribution of _trusted_judgments:")
9 df['_trusted_judgments'].hist()
10 plt.xlabel('_trusted_judgments')
11 plt.ylabel('Frequency')
12 plt.title('Distribution of _trusted_judgments')
13 plt.show()
```



12. How many reviews were finalized (_unit_state == 'finalized')?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 12. Number of Finalized Reviews (_unit_state == 'finalized')
8 finalized_reviews = df[df['_unit_state'] == 'finalized'].shape[0]
9 print("\n12. Number of finalized reviews:", finalized_reviews)
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
```

```
12. Number of finalized reviews: 1525
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

13. How many reviews have no policy violations?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 13. How Many Reviews Have No Policy Violations?
8 # Assuming reviews with empty string or NaN in 'policies_violated' have no violations
9 no_violations_count = df['policies_violated'].fillna('').str.strip().eq('').sum()
10 print(f"\n13. Number of reviews with no policy violations: {no_violations_count}")
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
```

```
13. Number of reviews with no policy violations: 0
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

14. Which policy is most commonly violated in each city?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 14. Most Commonly Violated Policy in Each City
8 def most_common_policy(violations):
9     return violations.str.split(',').explode().value_counts().idxmax()
10
11 common_violations_per_city = df.groupby('city')['policies_violated'].apply(most_common_policy)
12 print("\n14. Most commonly violated policy per city:")
13 print(common_violations_per_city)
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
```

```
14. Most commonly violated policy per city:
city
Atlanta      na
Chicago      na
Cleveland    SlowService
Dallas       na
Houston      na
Las Vegas    na
Los Angeles  na
New York     na
Portland     na
Name: policies_violated, dtype: object
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

15. What is the average confidence score for each type of violation?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 15. Average Confidence Score for Each Type of Violation
8 # Explode the policy violations and confidence scores into long format for averaging
9 df_long = df[['policies_violated', 'policies_violated:confidence']].copy()
10 df_long = df_long.dropna(subset=['policies_violated', 'policies_violated:confidence'])
11 df_long['policies_violated'] = df_long['policies_violated'].str.split(',')
12 df_long = df_long.explode('policies_violated')
13 df_long['policies_violated:confidence'] = pd.to_numeric(df_long['policies_violated:confidence'], errors='coerce')
14 average_confidence = df_long.groupby('policies_violated')['policies_violated:confidence'].mean()
15 print("\n15. Average confidence score for each type of violation:")
16 print(average_confidence)
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "C:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"

15. Average confidence score for each type of violation:
policies_violated
BadFood          0.876383
BadFood\rCost    NaN
BadFood\rFilthy  NaN
BadFood\rFilthy\rRudeService NaN
BadFood\rMissingFood NaN
...
SlowService\rScaryMcDs NaN
SlowService\rScaryMcDs\rBadFood NaN
SlowService\rScaryMcDs\rFilthy NaN
SlowService\rScaryMcDs\rRudeService NaN
na              0.909448
Name: policies_violated:confidence, Length: 140, dtype: float64
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```


16. Find the longest review by word count.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 16. Longest Review by Word Count
8 longest_review_index = df['review'].dropna().str.split().str.len().idxmax()
9 longest_review = df.loc[longest_review_index, 'review']
10 longest_review_length = len(longest_review.split())
11 print(f"\n16. Longest review by word count ({longest_review_length} words):")
12 print(longest_review)
```

PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"

16. Longest review by word count (914 words):

It's 2 for 1 review day! I can't review the Walmart without reviewing the McDonald's that resides inside. And boy is it some tantalizing stuff. I try not to give too many 1-star reviews but this McDonald's so deserves it. I hardly frequent a McDonald's at all (although this is my second McDonald's review in less than a week) but I need to let the peeps in on this one. No Yelper should go unwarned. Of course I came to McDonald's since I have those coupons for their new contraption, the McWrap. Which when I pronounce very fast and in my Scottish accent sounds much like 'McCrap'. After my second go round with this thing, I have made up my mind it is decidedly so. I need to write McDonald's execs to change the name on this thing immediately. My coupon for the day was buy a medium fries and a medium drink, get a McWrap of your choice for FREE. I have a hard time resisting free. I'm here to tell you...WORST ALMOST FREE MEAL OF MY LIFE !!! Seriously. It was the tail end of lunchtime so the line was still a bit busy. I patiently waited for almost 10 minutes...which is like forever in McDonald's time. At least two folks left the line while I was there. I know it's not McDonald's fault, but the couple in front of me took up majority of the line space because they were in motorized wheelchairs. And they were very loud. This doesn't bother me much. But...they were also very stinky. Yes, stinky. I have no problem with stinky in most cases (I occasionally help the homeless so I know) but the guy in front of me hadn't washed in DAYS if I were to guess. To each his own, and I know it's almost summer so I know we all want to save water...but DAMN. If I get a whiff and I have to take a step back? You might be S-T-I-N-K-Y. But who am I to judge? I took a big step back and waited my turn. Free is the motivation here. I finally get up there to the McDonald's cashier and present my coupon to the guy. He takes a look at it and says 'We don't usually take the promo coupons at this location but we'll let you use it.' Then he proceeds to put the coupon onto a stack of...more coupons! Mind boggling to say the least. He takes my order of a medium fries and medium drink, asks me what type of McWrap I want (Grilled Chicken & Bacon) and rings me up. \$2.91 later I'm given my receipt and asked to wait for my number to be called. I wait. And wait. #1! #2! #3! #5! #6! #7! #8! Obviously I'm #4 and 10 minutes later I am still waiting. Honestly, what the Hell are you guys doing back there? Mixing the masa for my fresh pressed tortilla? Feathering a chicken? Perhaps ripening a tomato? I ask the McDonald's runner and he tells me, "We're waiting for your McWrap." Of course. At this point I'm almost running late to be back at work, so the minute I get my order I bolt back to the office and have lunch at my desk. Thankfully I'm only three minutes away. I really, really should have just gone to the vending machine or eaten my stash of Reese's chocolate peanut butter eggs from last Easter. Talk about disappointment! I take out my McWrap from my bag and lay down the box of fries on my desk. I continue working while munching away on what is supposed to be McDonald's claim to fame. Their golden fries. I bite into one, eh a bit stale. I bit into another one, eh not enough salt. I bite into yet another one, eh cold?? I look down, and whaddaya know? My damn fry is not even barely cooked! White like Olivia Wilde in the winter. I look through the entire box and find at least half a dozen fries not cooked properly! Not golden brown! Not delicious! Not a

17. How many reviews mention "clean" or "dirty" in the text?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 17. Number of Reviews Mentioning "clean" or "dirty"
8 mentions_clean_dirty = df['review'].dropna().str.contains('clean|dirty', case=False).sum()
9 print(f"\n17. Number of reviews mentioning 'clean' or 'dirty': {mentions_clean_dirty}")
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"

17. Number of reviews mentioning 'clean' or 'dirty': 185
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

18. Are there duplicate reviews in the dataset?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 18. Count of Duplicate Reviews
8 duplicates = df.duplicated().sum()
9 print(f"\n18. Number of duplicate reviews: {duplicates}")
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "C:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
```

```
18. Number of duplicate reviews: 0
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

19. What is the median review length in characters?

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 19. Median Review Length in Characters
8 median_length = df['review'].dropna().str.len().median()
9 print(f"\n19. Median review length (characters): {median_length}")
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
```

```
19. Median review length (characters): 388.0
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh>
```

20. Plot the top 5 cities with the highest average number of violations per review.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 filename = 'McDonalds-Yelp-Sentiment-DFE.csv'
5 df = pd.read_csv(filename, encoding='ISO-8859-1')
6
7 # 20. Plot Top 5 Cities with Highest Average Number of Violations per Review
8 avg_violations_per_city = df.groupby('city')['policies_violated'].apply(
9     lambda x: x.str.split(',').apply(len).mean()
10 ).nlargest(5)
11
12 print("\n20. Top 5 cities with highest average violations per review:")
13 print(avg_violations_per_city)
14
15 avg_violations_per_city.plot(kind='bar', color='skyblue')
16 plt.xlabel('City')
17 plt.ylabel('Average Number of Violations per Review')
18 plt.title('Top 5 Cities with Highest Average Number of Violations per Review')
19 plt.xticks(rotation=45)
20 plt.tight_layout()
21 plt.show()
```

```
PS C:\Users\swapn\OneDrive\Desktop\EDS Theory Activity No.1\Adarsh> & C:/Users/swapn/AppData/Local/Programs/Python/Python313/python.exe "c:/Users/swapn/OneDrive/Desktop/EDS Theory Activity No.1/Adarsh/Final.py"
```

```
20. Top 5 cities with highest average violations per review:
city
Atlanta      1.0
Chicago      1.0
Cleveland    1.0
Dallas       1.0
Houston      1.0
Name: policies_violated, dtype: float64
```

