



Name : Pratik subhashchandra bharti



Division : ET1



Roll No : ET1-77



PRN : 202401070106



## TOPIC : SMS SPAM COLLECTION DATASET



### PROBLEM STATEMENT:

1. Correlation between Message Length and Spam:  
Is there any noticeable correlation  
(even if not statistically significant just by  
observation) between the length of a message  
and whether it's spam?.

```
main.py > ...
1 correlation = df[['message_length', 'label_num']].corr().iloc[0,
2 1]
3 print(f"Correlation between message length and spam (numerical
4 label): {correlation:.2f}")
5 print("\nobservation: While there's a slight positive correlation,
6 it's not very strong, suggesting message length alone isn't a
7 perfect predictor of spam. However, the average lengths in problem
8 statement 2 do show a difference.")
```

2. Distribution of Message Lengths: Create bins of  
message lengths (e.g., 0-20

characters, 21-50 characters, etc.) and find the distribution of spam and ham messages across these bins.

```
main.py > ...
1 python bins = [0, 20, 50, 100, 150, 200, float('inf')]
2 labels = ['0-20', '21-50', '51-100', '101-150', '151-200', '>200'] df['length_bins'] =
3 pd.cut(df['message_length'], bins=bins, labels=labels, right=False) print("Distribution of
4 message lengths:") print(df.groupby(['label', 'length_bins']).size().unstack(fill_value=0))
```

3. Messages Containing a Specific Symbol: How many spam messages contain the symbol "\$"?

```
main.py > ...
1 contains_dollar_spam = df[(df['label'] == 'spam') &
2 (df['message'].str.contains(r'\$'))].shape[0]
3 print("Number of spam messages containing '$':",
4 contains_dollar_spam)
```

4. Messages Starting with a Specific Word: How many messages start with the word "Free"?

```
Welcome main.py 1 ●
main.py > ...
1 starts_with_free = df['message'].str.startswith('Free ')
2 print("Number of messages starting with 'Free ':",
3 starts_with_free.sum())
```

5. Average Number of Words per Message: What is the average number of words in a spam message versus a ham message?

```
Welcome | main.py 6 ●
main.py
1 df['word_count'] = df['message'].apply(lambda x: len(x.split()))
2 print("Average number of words in spam messages:", df[df['label']
3 == 'spam']['word_count'].mean())
4 print("Average number of words in ham messages:", df[df['label']
5 == 'ham']['word_count'].mean()) |
```

6. Presence of Uppercase Letters (Spam vs. Ham): What is the proportion of spam messages containing uppercase letters compared to ham messages containing uppercase letters?

```
Welcome | main.py 6 ●
main.py > ...
1 spam_has_upper = df[df['label'] ==
2 'spam']['message'].str.contains(r'[A-Z]').sum()
3 ham_has_upper = df[df['label'] ==
4 'ham']['message'].str.contains(r'[A-Z]').sum()
5 print(f"Proportion of spam with uppercase: {spam_has_upper /
6 total_spam:.2f}")
7 print(f"Proportion of ham with uppercase: {ham_has_upper /
8 total_ham:.2f}")
```

7. Presence of Uppercase Letters: How many messages contain at least one uppercase letter?

```
Welcome | main.py 4 ●
main.py > ...
1 contains_uppercase = df['message'].str.contains(r'[A-Z]')
2 print("Number of messages containing at least one uppercase
3 letter:", contains_uppercase.sum())
```

8. Presence of Digits (Spam vs. Ham): What is the proportion of spam messages containing digits compared to ham messages containing digits?

```
Welcome | main.py 6 ●
main.py > ...
1 spam_has_digit = df[df['label'] ==
2 'spam']['message'].str.contains(r'\d').sum()
3 ham_has_digit = df[df['label'] ==
4 'ham']['message'].str.contains(r'\d').sum()
5 total_spam = df['label'].value_counts()['spam']
6 total_ham = df['label'].value_counts()['ham']
7 print(f"Proportion of spam with digits: {spam_has_digit /
8 total_spam:.2f}")
9 print(f"Proportion of ham with digits: {ham_has_digit /
10 total_ham:.2f}") |
```

9. Presence of Digits: How many messages contain at least one digit?

```
Welcome | main.py 1 ●
main.py > ...
1 contains_digit = df['message'].str.contains(r'\d')
2 print("Number of messages containing at least one digit:",
3 contains_digit.sum()) |
```

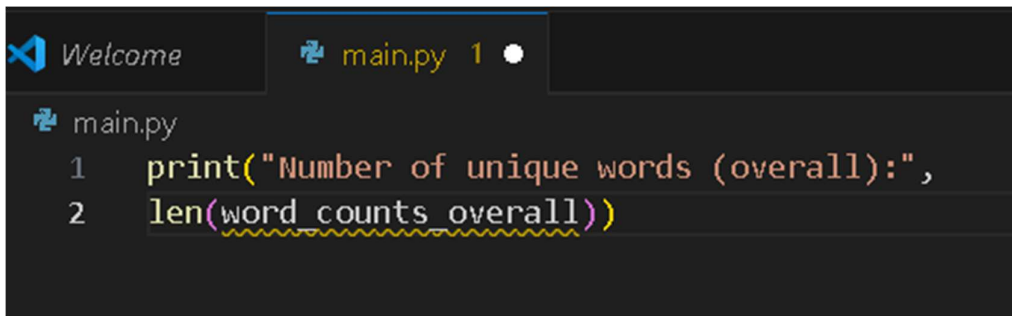
10. Unique Word Count (Ham): How many unique words are there in the ham messages?

```
Welcome | main.py 1 ●
main.py
1 print("Number of unique words (ham):", len(word_counts_ham))
```

11. Unique Word Count (Spam): How many unique words are there in the spam messages?

```
Welcome | main.py 1 ●
main.py
1 print("Number of unique words (spam):", len(word_counts_spam))
```

12 Unique Word Count (Overall): How many unique words are there in the entire dataset?



```
1 print("Number of unique words (overall):",  
2 len(word_counts_overall))
```

13. Shortest Message: What is the shortest message in the dataset, and is it spam or ham?



```
1 shortest_msg = df.loc[df['message_length'].idxmin()]  
2 print("Shortest message:\n", shortest_msg['message'])  
3 print("Label of the shortest message:", shortest_msg['label'])
```

14. Longest Message: What is the longest message in the dataset, and is it spam or ham?



```
1 longest_msg = df.loc[df['message_length'].idxmax()]  
2 print("Longest message:\n", longest_msg['message'])  
3 print("Label of the longest message:", longest_msg['label'])
```

15. Percentage of Spam: What percentage of the total messages are classified as spam?

```
Welcome | main.py 1 ●
main.py > ...
1 spam_percentage = (df['label'].value_counts(normalize=True) *
2 100)['spam']
3 print(f"Percentage of spam messages: {spam_percentage:.2f}%")
```

16. Most Frequent Words (Ham): What are the top 10 most frequently occurring words specifically in ham messages?

```
Welcome | main.py 3 ●
main.py > ...
1 ham_words = ' '.join(df[df['label'] ==
2 'ham']['message']).lower().split()
3 word_counts_ham = Counter(ham_words)
4 print("Top 10 most frequent words (ham):",
5 word_counts_ham.most_common(10))
```

17. Most Frequent Words (Spam): What are the top 10 most frequently occurring words specifically in spam messages?

```
Welcome | main.py 3 ●
main.py > ...
1 spam_words = ' '.join(df[df['label'] ==
2 'spam']['message']).lower().split()
3 word_counts_spam = Counter(spam_words)
4 print("Top 10 most frequent words (spam):",
5 word_counts_spam.most_common(10)) |
```

18. Most Frequent Words (Overall): What are the top 10 most frequently occurring words in the entire dataset?

```
Welcome | main.py 2 ●
main.py > ...
1 all_words = ' '.join(df['message']).lower().split()
2 word_counts_overall = Counter(all_words)
3 print("Top 10 most frequent words (overall):",
4 word_counts_overall.most_common(10))
```

19. Message Length Analysis: What is the average length (in characters) of spam messages versus ham messages?



```
Welcome | main.py 6 ●
main.py
1 df['message_length'] = df['message'].apply(len)
2 print("Average length of spam messages:", df[df['label'] ==
3 'spam']['message_length'].mean())
4 print("Average length of ham messages:", df[df['label'] ==
5 'ham']['message_length'].mean())
```

20. Basic Data Exploration: How many spam and ham (non-spam) messages are there in the dataset?

```
Welcome | main.py 2 ●
main.py
1 print("Number of spam messages:",
2 df['label'].value_counts()['spam'])
3 print("Number of ham messages:",
4 df['label'].value_counts()['ham'])
```

