

Table of Contents



Problem Overview

Describing our goal for performing NLP analysis



Datasets/Visualizations

Cleaning data, Word Clouds & Frequencies, Scatterplots



Sentiment Analysis

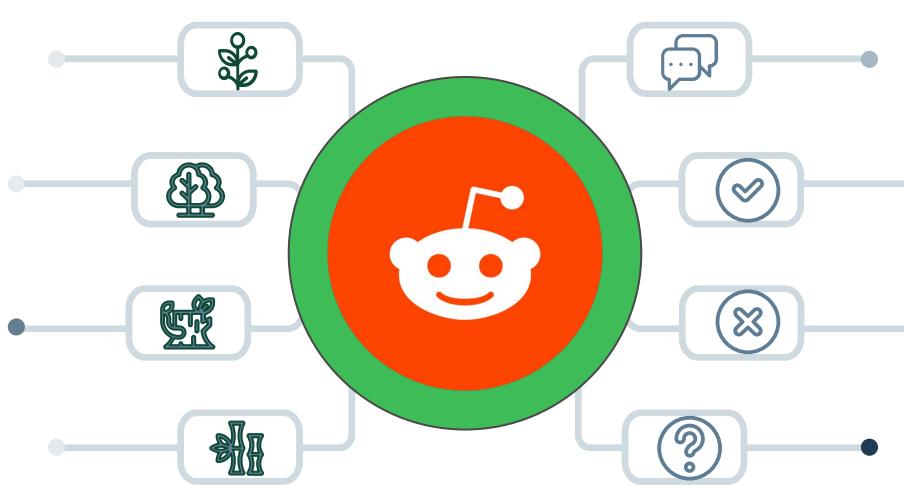
VADER process



Conclusion

Results, insights, challenges





Problem Overview

Determine the sentiment of Reddit comments regarding environmental policies on three different subreddits using esg_scores and the averages of the sentiment compound scores

...

Datasets

Environment

Article comments from environment subreddit

Climate

Article comments from climate subreddit

Sustainability

Article comments from sustainability subreddit

We used a combination of two separate datasets we found on Kaggle:

1. "🌱 Greenwashing on Reddit 🌎"
 - Greenwashing Articles and Reddit's Posts and Comments up to 03/2025
2. "S&P 500 Firms ESG Sustainability Reports Dataset" (preview below)



Cleaning Datasets

Removing stop words

- a. Stop words csv
- b. Adding custom stop words

stopwords.txt

```
1 x
2 y
3 your
4 yours
5 yourself
6 yourselves
7 you
8 yond
9 yonder
10 yon
11 ye
12 yet
13 z
14 zillion
15 j
16 u
17 umpteen
18 usually
19 us
20 username
21 uponed
22 upons
23 unoning
```

```
1 # snp_df Word Cloud
2
3 from wordcloud import WordCloud, STOPWORDS
4 import matplotlib.pyplot as plt
5 import pandas as pd
6
7 # Environment Subreddit Articles Comments Word Cloud
8
9 comments_df = pd.read_csv('/work/reddit_folders/environmentSubreddit/Comments/environmentArticlesComments.csv')
10 comments = ' '.join(comments_df['Text'].astype(str))
11
12 #stopwords
13 stopwords_df = pd.read_csv('/work/stopwords.txt', header = None)
14 custom_stopwords = STOPWORDS.union(set(stopwords_df[0].astype(str).str.lower()))
15 commentscloud = WordCloud(
16     width=800,
17     height=400,
18     background_color='white',
19     colormap='viridis',
20     max_words=100,
21     stopwords=custom_stopwords
22 ).generate(comments)
23 plt.figure(figsize=(15, 8))
24 plt.imshow(commentscloud, interpolation='bilinear')
25 plt.axis('off')
26 plt.title('Word Cloud - Environment Subreddit Articles Comments', fontsize=20)
27 plt.tight_layout()
28 plt.show()
```

Visualizations

- a. Word Clouds 
- b. Word Frequencies 
- c. Scatter Plots on relationship between ESG scores 





Word Clouds

- a. Environment Subreddit Comment Word Cloud
- b. Climate Subreddit Comment Word Cloud
- c. Sustainability Subreddit Comment Word Cloud
- d. S&P 500 ESG Corporate Reports Word Cloud

Goal: Visually compare how the public talks about sustainability versus how companies officially describe them.

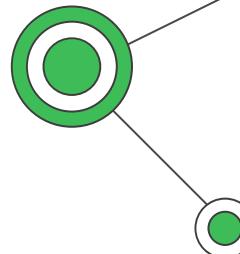


Environment Subreddit Comment Word Cloud

```
1 # snp_df Word Cloud
2
3 from wordcloud import WordCloud, STOPWORDS
4 import matplotlib.pyplot as plt
5 import pandas as pd
6
7 # Environment Subreddit Articles Comments Word Cloud
8
9 comments_df = pd.read_csv('/work/reddit_folders/environmentSubreddit/Comments/environmentArticlesComments.csv')
10 comments = ' '.join(comments_df['Text'].astype(str))
11
12 #stopwords
13 stopwords_df = pd.read_csv('/work/stopwords.txt', header = None)
14 custom_stopwords = STOPWORDS.union(set(stopwords_df[0].astype(str).str.lower()))
15 commentscloud = WordCloud(
16     width=800,
17     height=400,
18     background_color='white',
19     colormap='viridis',
20     max_words=100,
21     stopwords=custom_stopwords
22 ).generate(comments)
23 plt.figure(figsize=(15, 8))
24 plt.imshow(commentscloud, interpolation='bilinear')
25 plt.axis('off')
26 plt.title('Word Cloud - Environment Subreddit Articles Comments', fontsize=20)
27 plt.tight_layout()
28 plt.show()
```



Environment Subreddit Comment Word Cloud



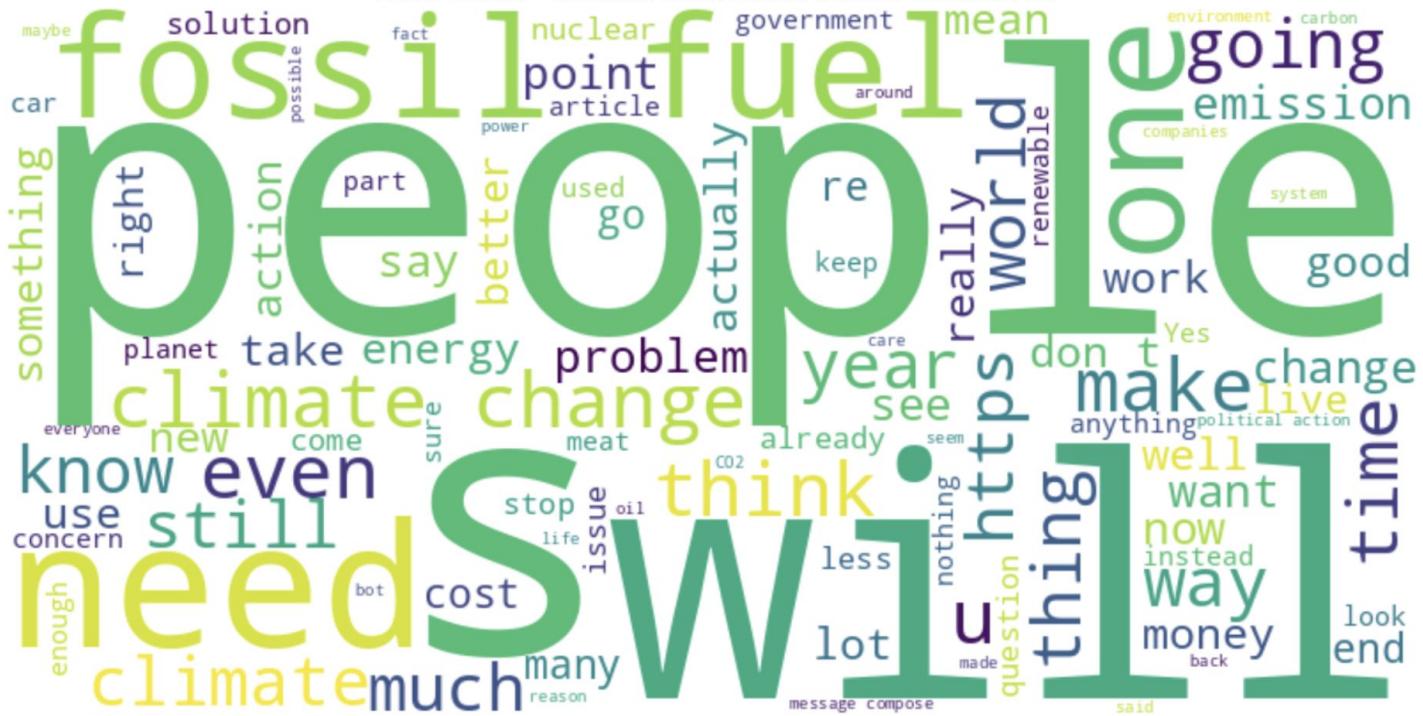
Climate Subreddit Comment Word Cloud

```
1 # Climate Subreddit Articles Comments Word Cloud
2
3 climateComments_df = pd.read_csv('/work/reddit_folders/climateSubreddit/Comments/climateArticlesComments.csv')
4 climateComments = ' '.join(climateComments_df['Text'].astype(str))
5 climateCommentsCloud = WordCloud(width=800, height=400,
6                                 background_color='white',
7                                 colormap='viridis',
8                                 max_words=100).generate(climateComments)
9 plt.figure(figsize=(15, 8))
10 plt.imshow(climateCommentsCloud, interpolation='bilinear')
11 plt.axis('off')
12 plt.title('Word Cloud - Climate Subreddit Articles Comments', fontsize=20)
13 plt.tight_layout()
14 plt.show()
```



Climate Subreddit Comment Word Cloud

Word Cloud - Climate Subreddit Articles Comments





Sustainability Subreddit Comment Word Cloud

```
1 # Sustainability Subreddit Articles Comments Word Cloud
2
3 sustainabilityComments_df = pd.read_csv('/work/reddit_folders/sustainabilitySubreddit/Comments/sustainabilityArticlesComments_df.csv')
4 sustainabilityComments = ' '.join(climateComments_df['Text'].astype(str))
5 sustainabilityCommentsCloud = WordCloud(width=800, height=400,
6                                         background_color='white',
7                                         colormap='viridis',
8                                         max_words=100).generate(sustainabilityComments)
9 plt.figure(figsize=(15, 8))
10 plt.imshow(sustainabilityCommentsCloud, interpolation='bilinear')
11 plt.axis('off')
12 plt.title('Word Cloud - Sustainability Subreddit Articles Comments', fontsize=20)
13 plt.tight_layout()
14 plt.show()
```



Sustainability Subreddit Comment Word Cloud





S&P Corporate Reports Word Cloud

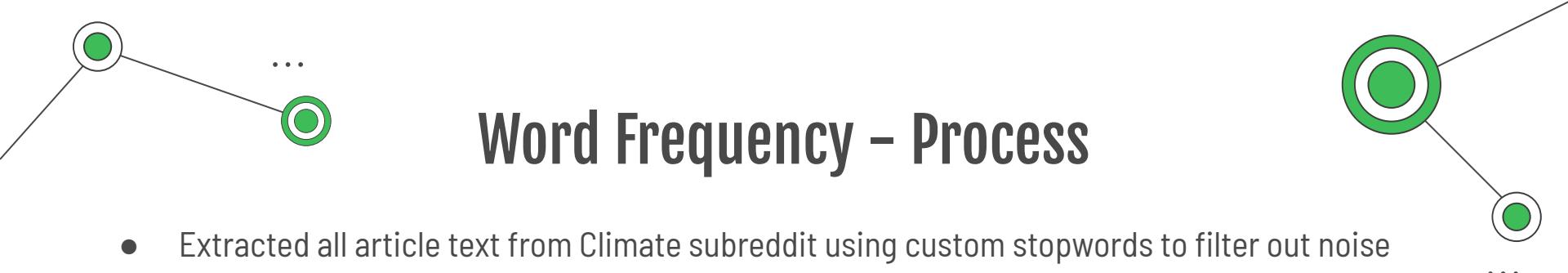
```
1 # Combine all text from the column into one string
2 text = ' '.join(snp_df['preprocessed_content'].astype(str))
3
4
5 # Create the word cloud
6 wordcloud = WordCloud(stopwords=custom_stopwords,
7                         width=800, height=400,
8                         background_color='white',
9                         colormap='viridis',
10                        max_words=100).generate(text)
11
12 # Display the word cloud
13 plt.figure(figsize=(15, 8))
14 plt.imshow(wordcloud, interpolation='bilinear')
15 plt.axis('off')
16 plt.title('Word Cloud', fontsize=20)
17 plt.tight_layout()
18 plt.show()
```



Word Frequency – Process

- Extracted all article text from Climate subreddit using custom stopwords to filter out noise

```
1 climate_change = pd.read_csv('reddit_folders/climatechangeSubreddit/Articles/climatechangeArticles.csv')
2 climate_change.head()
```

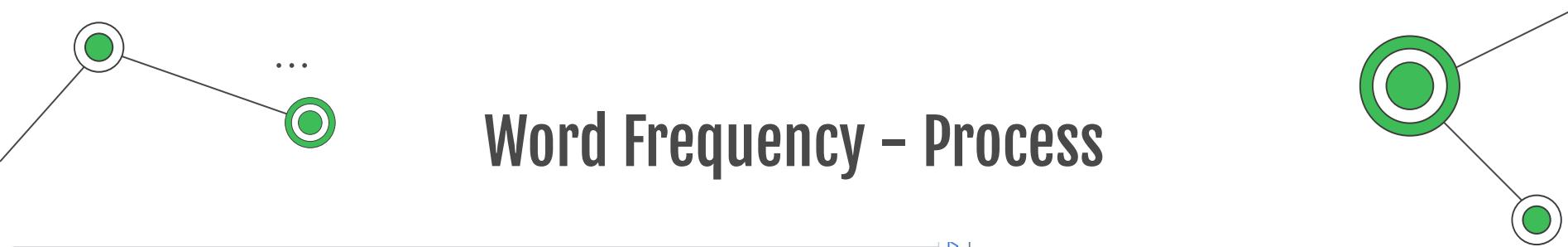


	Article_Id int64	Article_Url object	Post_id object	Article_Title object	Article_Content o...	
0	1	https://www.eur...	e74s78	'Do no harm': Nuc...	Euractiv +For indi...	
1	2	https://insidecli...	e83dht	Lawyers Challeng...	Lawyers Challeng...	
2	3	https://www.cbc....	ehe1pu	New industry dev...	When search sug...	
3	4	https://medium.c...	gh21y2	Microsoft's 'Carb...	Microsoft's 'Carb...	
4	5	https://2030.bui...	jvsc81	5 Ways to Avoid ...	Greenwashing: W...	

5 rows, 5 cols 10 / page << < Page 1 of 1 > >> ⚡ Form

```
1 climate_change["Article_Content"][0]
```

```
"Euractiv +For individuals Euractiv ProFor corporations Looking to access paid articles across multiple policy topics? Int
```

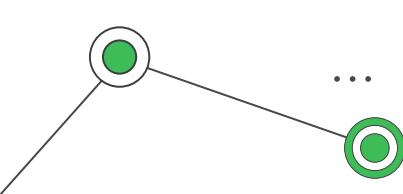


Word Frequency – Process

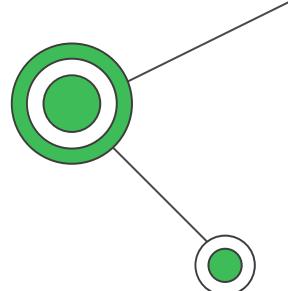
```
1 STOPWORDS = {  
2     'the', 'a', 'an', 'and', 'or', 'to', 'of', 'in', 'is', 'it', 'for', 'on', 'that', 'this', 'with', 'as',  
3     'are', 'was', 'be', 'by', 'from', 'at', 'if', 'but', 'so', 'not', 'you', 'i', 'we', 'they', 'he', 'she'  
4 }  
5  
6 word_counts = {}  
7  
8 for article in climate_change["Article_Content"]:  
9     words = re.findall(r"[a-z]+", article.lower())  
10    #words = article.lower().split() # lowercase + split on whitespace  
11    for word in words:  
12        if word in STOPWORDS:  
13            continue  
14        word_counts[word] = word_counts.get(word, 0) + 1  
15  
16 # This sorts our dictionary based off the corresponding frequency, and turns it into a list.  
17 word_counts_show = sorted(word_counts.items(), key=lambda item: item[1], reverse = True)  
18 word_counts_show[:10] # Display the top 20 most frequent words
```

```
[('s', 589),  
 ('climate', 399),  
 ('emissions', 335),  
 ('have', 282),  
 ('more', 269),  
 ('carbon', 266),  
 ('co', 258),  
 ('has', 227),  
 ('energy', 224),  
 ('can', 209)]
```

-
- ```
graph LR; A(()); B(()); C(());
```
- Counted word frequencies by
    - Extract alphabetical tokens
    - Remove stopwords
    - Aggregate their occurrence across all articles



# Word Frequency – Process



```
1 def wordDictionary(df, columnName):
2 STOPWORDS = {'the', 'a', 'an', 'and', 'or', 'to', 'of', 'in', 'is', 'it', 'for', 'on', 'that', 'this', 'with', 'as',
3 'are', 'was', 'be', 'by', 'from', 'at', 'if', 'but', 'so', 'not', 'you', 'i', 'we', 'they', 'he', 'she'}
4 word_counts = {}
5 for article in df[columnName]:
6 words = re.findall(r"[a-z]+", article.lower())
7 #words = article.lower().split() # lowercase + split on whitespace
8 for word in words:
9 if word in STOPWORDS:
10 continue
11 word_counts[word] = word_counts.get(word, 0) + 1
12
13 return word_counts
```

- Extract tokens, remove stopwords, and compute word frequencies

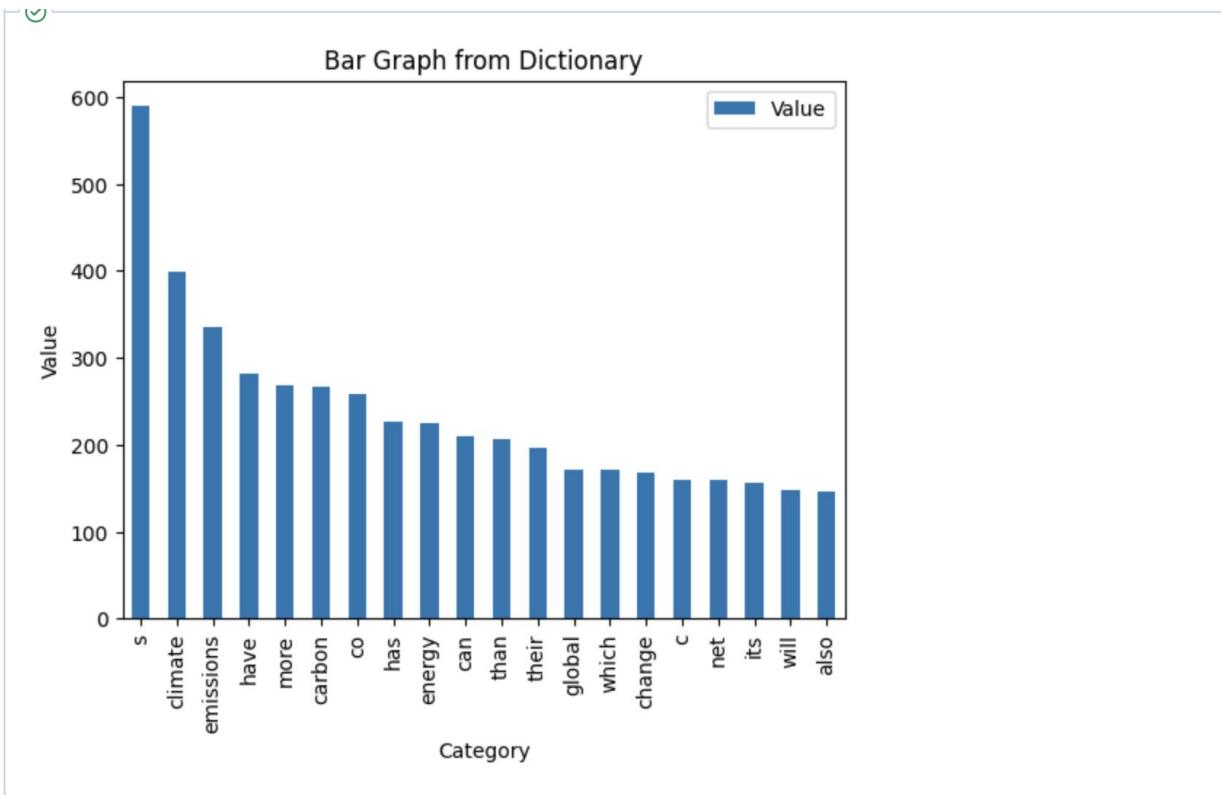
```
1 df = pd.DataFrame.from_dict(word_counts, orient='index', columns=['Value'])
2 df = df.sort_values(by=['Value'], ascending=False)
3 df.head(5)
4 df = df[:20]
5
```

- Convert the frequency into a dataframe

```
1 df.plot(kind='bar')
2 plt.title('Bar Graph from Dictionary')
3 plt.xlabel('Category')
4 plt.ylabel('Value')
5 plt.show()
```

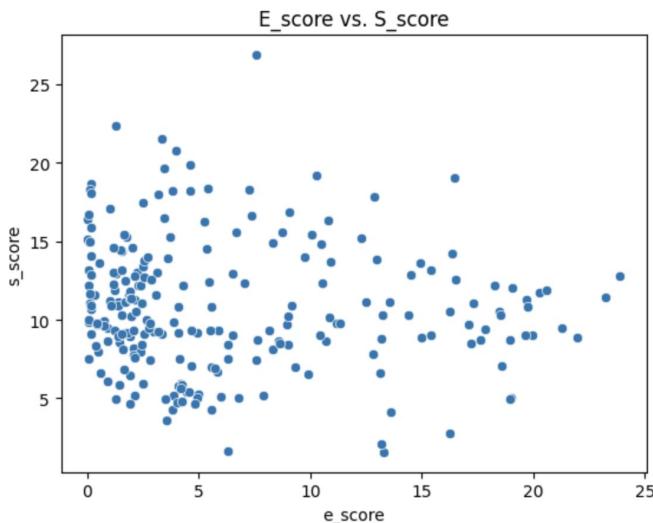
- Visualize the top word frequency with a bar chart

# Word Frequency - Result

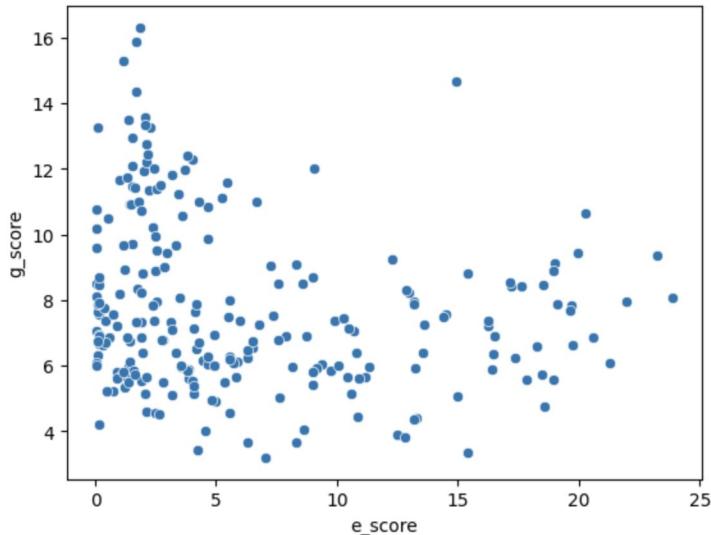


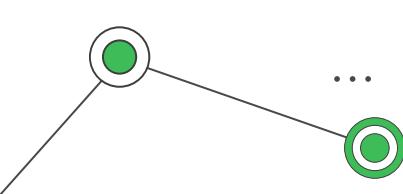
# Scatter Plots – Relationship between ESG Scores

```
1 import seaborn as sns
2
3 sns.scatterplot(x="e_score", y="s_score", data=snp_grouped)
4 plt.title("E_score vs. S_score")
5 plt.show()
```

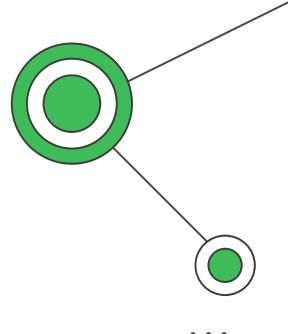


```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 sns.scatterplot(x="e_score", y="g_score", data=snp_grouped)
5 plt.show()
```



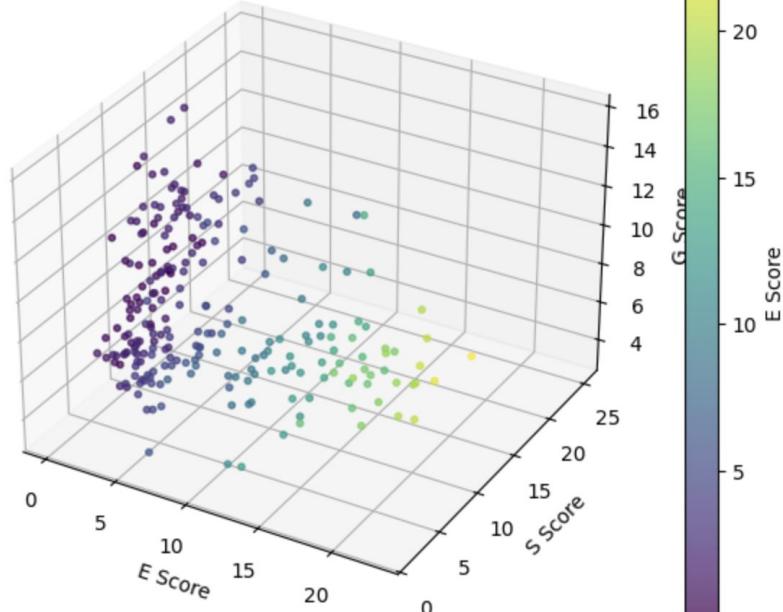


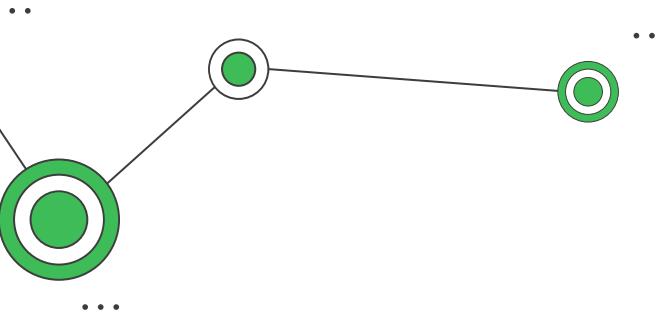
# 3D Scatter Plot



```
1 from mpl_toolkits.mplot3d import Axes3D
2
3 fig = plt.figure(figsize=(8, 6))
4 ax = fig.add_subplot(111, projection='3d')
5
6 sc = ax.scatter(
7 avg['e_score'],
8 avg['s_score'],
9 avg['g_score'],
10 c=avg['e_score'],
11 cmap='viridis',
12 s = 10, alpha = 0.7
13)
14 ax.set_xlabel('E Score')
15 ax.set_ylabel('S Score')
16 ax.set_zlabel('G Score')
17 ax.set_title('3D Scatter Plot of ESG Scores')
18 fig.colorbar(sc, ax=ax, label = 'E Score')
19 plt.show()
```

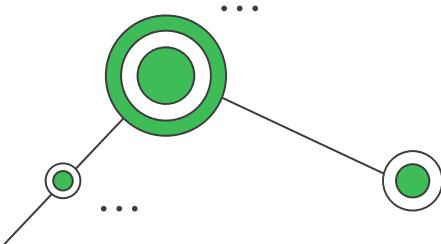
3D Scatter Plot of ESG Scores

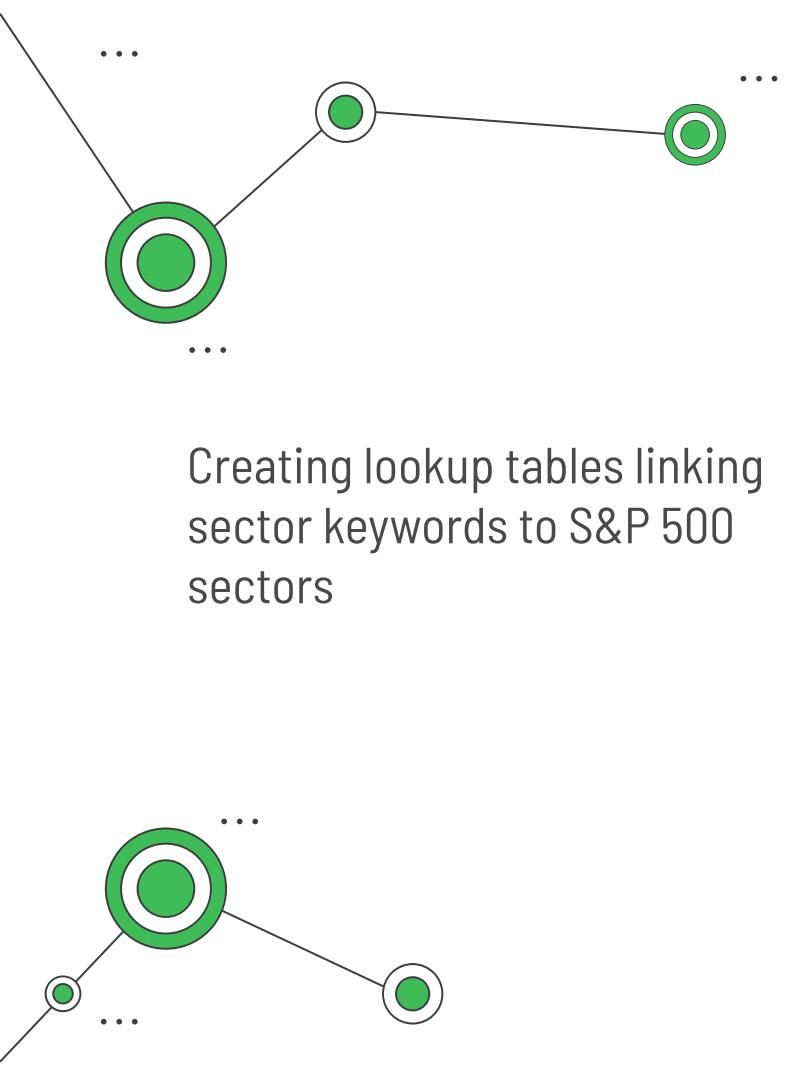




## Categorizing Reddit comments into sectors 🍰

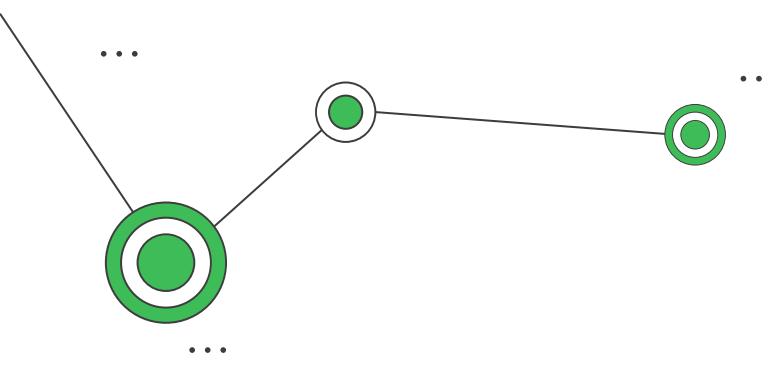
- To gain structured insights
- To learn what industries people care about the most





## Creating lookup tables linking sector keywords to S&P 500 sectors

```
1 def extract_sector_info(sp500_df):
2 """
3 Build lookups for sector keywords and map companies to sectors.
4 Similar to extract_company_info but for sectors.
5
6 Returns:
7 - sector_keywords: dict mapping keywords to sectors
8 - sector_lookup: dict mapping tickers to sectors (for reference)
9 """
10
11 # Define comprehensive sector-specific keywords
12 sector_keyword_map = {
13 'Technology': [
14 'tech', 'software', 'cloud', 'ai', 'artificial intelligence',
15 'semiconductor', 'chip', 'chips', 'data center', 'cybersecurity', 'saas',
16 'hardware', 'computing', 'digital', 'platform', 'app', 'apps',
17 'mobile', 'internet', 'web', 'coding', 'programming'
18],
19 'Health Care': [
20 'pharma', 'pharmaceutical', 'biotech', 'healthcare', 'health care',
21 'medical', 'drug', 'drugs', 'medicine', 'hospital', 'clinical',
22 'vaccine', 'vaccines', 'therapy', 'patient', 'device', 'diagnostic',
23 'treatment', 'doctor', 'healthcare provider'
24],
25 'Financials': [
26 'bank', 'banking', 'financial', 'finance', 'insurance', 'credit',
27 'loan', 'loans', 'mortgage', 'investment', 'asset management',
28 'wealth', 'broker', 'brokerage', 'lending', 'fintech',
29 'payment', 'payments', 'credit card'
30],
31 'Consumer Discretionary': [
32 'retail', 'ecommerce', 'e-commerce', 'consumer', 'shopping', 'store',
33 'restaurant', 'restaurants', 'auto', 'automotive', 'car', 'cars',
34 'hotel', 'hotels', 'travel', 'entertainment', 'luxury', 'apparel',
35 'clothing', 'fashion', 'furniture', 'home improvement'
36],
37 'Communication Services': [
38 'telecom', 'telecommunications', 'media', 'streaming', 'content',
39 'advertising', 'social media', 'network', 'broadcast', 'cable',
40 'wireless', '5g', 'internet service', 'video streaming'
41],
42 'Industrials': [
43 'industrial', 'manufacturing', 'aerospace', 'defense',
44 'construction', 'machinery', 'transportation', 'logistics',
45 'shipping', 'railroad', 'airline', 'airlines', 'freight'
46],
47 }
```



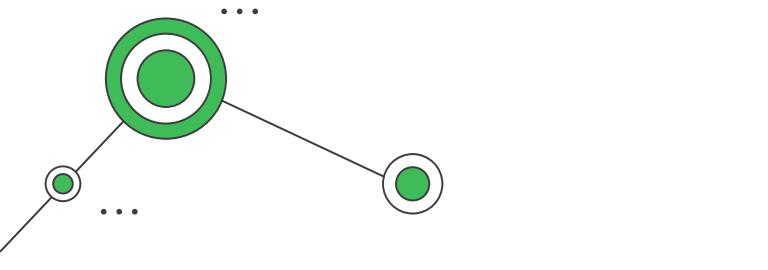
## Creating sector lookup dictionaries

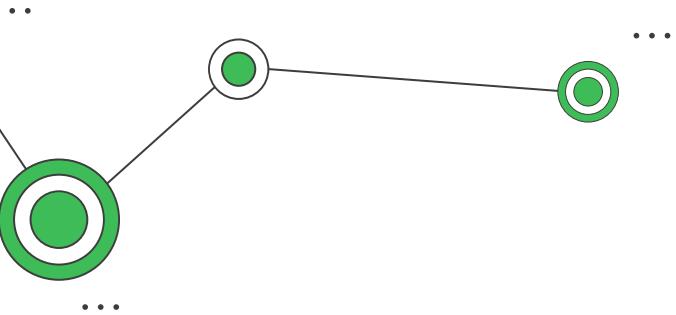
```
],
 'Consumer Staples': [
 'food', 'beverage', 'grocery', 'household', 'staples',
 'packaged goods', 'tobacco', 'personal care', 'consumer goods',
 'snacks', 'drinks'
],
 'Energy': [
 'oil', 'gas', 'energy', 'renewable', 'solar', 'wind',
 'petroleum', 'drilling', 'refining', 'electric utility',
 'fossil fuel', 'clean energy', 'crude'
],
 'Utilities': [
 'utility', 'utilities', 'electric', 'electricity', 'power',
 'water', 'gas utility', 'power grid'
],
 'Real Estate': [
 'real estate', 'reit', 'reits', 'property', 'properties',
 'commercial real estate', 'residential', 'mortgage reit',
 'landlord', 'housing'
],
 'Materials': [
 'materials', 'mining', 'chemical', 'chemicals', 'metals',
 'steel', 'commodity', 'commodities', 'packaging', 'paper',
 'construction materials', 'aluminum', 'copper'
]
}

Create reverse lookup: keyword -> sector
sector_keywords = {}
for sector, keywords in sector_keyword_map.items():
 for keyword in keywords:
 sector_keywords[keyword.lower()] = sector

Create ticker -> sector lookup (same as before)
sector_lookup = {}
for index, row in sp500_df.iterrows():
 ticker = str(row['ticker']).upper()
 sector = row['Sector'] if ('Sector' in row and pd.notna(row['Sector'])) else 'Unknown'
 sector_lookup[ticker] = sector

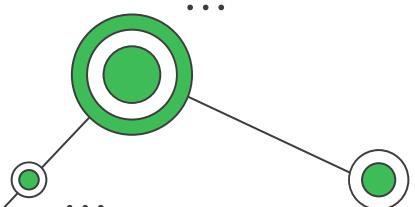
return sector_keywords, sector_lookup
```





## Identifying sector mentions and extracting context

Return mentions (at the end)



```
def find_sector_mentions_with_context(text, sector_keywords):
 """
 Find sector mentions and extract surrounding context.
 Similar to find_company_mentions_with_context but for sectors.

 Returns list of dicts with sector and its context window
 """
 if pd.isna(text):
 return []

 text = str(text)
 doc = nlp(text)
 mentions = []

 # Track which sectors we've already found to avoid duplicates
 found_sectors = {}

 # Search for each keyword in the text
 for keyword, sector in sector_keywords.items():
 pattern = r'\b' + re.escape(keyword) + r'\b'
 for match in re.finditer(pattern, text, re.IGNORECASE):
 # Skip if we already found this sector
 if sector in found_sectors:
 continue

 # Find the sentence containing this match
 for sent in doc.sents:
 if sent.start_char <= match.start() <= sent.end_char:
 # Get surrounding sentences for better context
 sent_index = list(doc.sents).index(sent)
 all_sents = list(doc.sents)

 context_start = max(0, sent_index - 1)
 context_end = min(len(all_sents), sent_index + 2)
 context = ' '.join([s.text for s in all_sents[context_start:context_end]])

 mentions.append({
 'sector': sector,
 'mention_text': match.group(),
 'sentence': sent.text,
 'context': context,
 'char_start': match.start(),
 'char_end': match.end()
 })
 found_sectors[sector] = True
 break
```

```

def group_comments_by_sector(comments_df, sp500_df, keep_single_sector_only=False):
 """
 Main function to group Reddit comments by sector.
 Similar to group_comments_by_company but for sectors.

 Parameters:
 - comments_df: DataFrame with Reddit comments
 - sp500_df: DataFrame with S&P 500 company data
 - keep_single_sector_only: If True, only keeps comments mentioning exactly 1 sector

 Returns:
 - enriched_comments: DataFrame with sector info added to each comment
 - sector_groups: Dictionary with sectors as keys, lists of comments as values
 """
 # Extract sector information
 sector_keywords, sector_lookup = extract_sector_info(sp500_df)

 # Extract sector mentions with context
 comments_df['sector_mentions'] = comments_df['Text'].apply(
 lambda x: find_sector_mentions_with_context(x, sector_keywords)
)

 # Expand to one row per sector mention
 expanded_rows = []

 for index, row in comments_df.iterrows():
 mentions = row['sector_mentions']

 if keep_single_sector_only and len(mentions) != 1:
 continue # Skip comments with 0 or multiple sectors

 if len(mentions) == 0:
 continue # Skip comments with no sector mentions

 # Create a row for each sector mentioned
 for mention in mentions:
 new_row = row.copy()
 new_row['mentioned_sector'] = mention['sector']
 new_row['sector_mention_text'] = mention['mention_text']
 new_row['sector_context'] = mention['context']
 new_row['sector_sentence'] = mention['sentence']
 new_row['is_multi_sector'] = len(mentions) > 1
 expanded_rows.append(new_row)

 expanded_df = pd.DataFrame(expanded_rows)

```

```

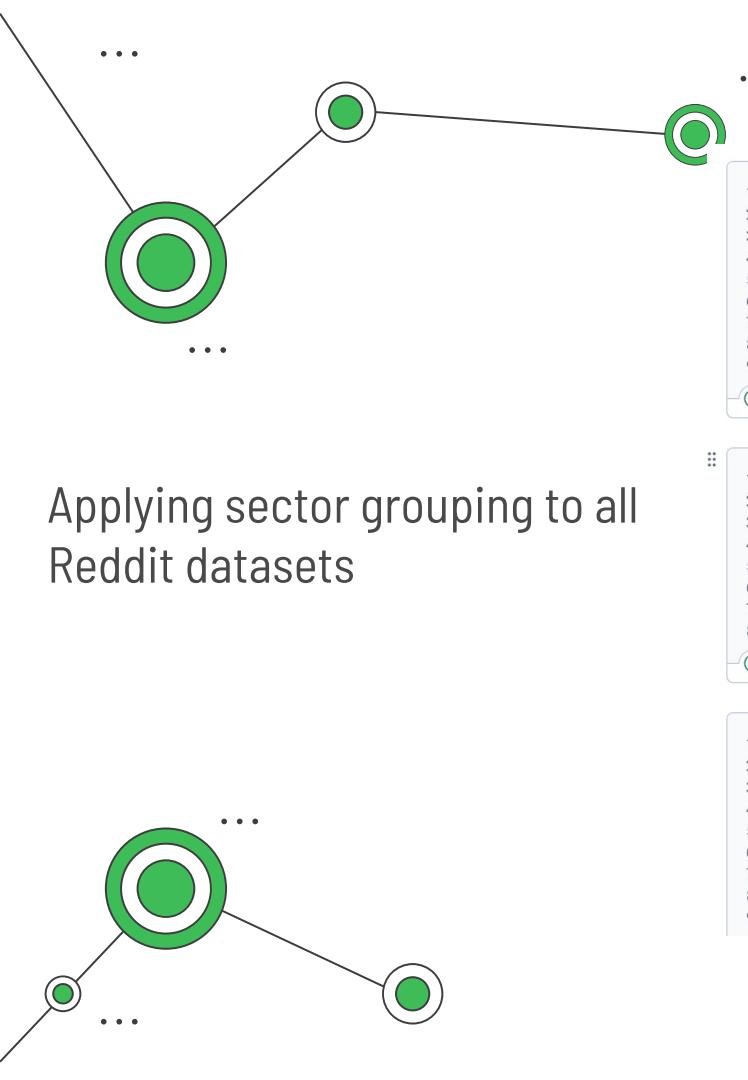
expanded_df = pd.DataFrame(expanded_rows)

Create grouped datasets by sector
sector_groups = defaultdict(list)
for index, row in expanded_df.iterrows():
 sector = row['mentioned_sector']
 sector_groups[sector].append({
 'comment_id': row.get('Parent_id', index),
 'text': row['Text'],
 'author': row.get('Author', ''),
 'score': row.get('Score', 0),
 'timestamp': row.get('Timestamp', ''),
 'sector': row['mentioned_sector'],
 'sector_context': row['sector_context'],
 'is_multi_sector': row['is_multi_sector']
 })
return expanded_df, sector_groups

def create_sector_summary(sector_groups):
 """
 Create a summary of comments by sector.
 Similar to create_company_summary but for sectors.
 """
 summary = []
 for sector, comments in sector_groups.items():
 summary.append({
 'Sector': sector,
 'Comment_Count': len(comments),
 'Single_Sector_Comments': sum(1 for comment in comments if not comment['is_multi_sector']),
 'Multi_Sector_Comments': sum(1 for comment in comments if comment['is_multi_sector']),
 'Total_Score': sum(comment['score'] for comment in comments),
 'Avg_Score': sum(comment['score'] for comment in comments) / len(comments) if comments else 0,
 'Unique_Authors': len(set(comment['author'] for comment in comments))
 })
 return pd.DataFrame(summary).sort_values('Comment_Count', ascending=False)

```

## Grouping comments by sector and create sector-level summaries



## Applying sector grouping to all Reddit datasets

```
1 #TODO call the function here on all of the dataframes we have, try to call it on the regular dataframe (so not the sam
2 #if it takes way too long then call the function on the sampled version
3 #make sure to store these new tables into different variables we can use later on
4 sp500_df = pd.read_csv("/work/snp_500.csv")
5 envCommentsDF_grouped = group_comments_by_sector(
6 envCommentsDF, sp500_df, keep_single_sector_only=False
7)[0]
8 envCommentsDF_grouped.to_csv("envCommentsDF_grouped.csv", index=False)
9
```

```
1 #TODO call the function here on all of the dataframes we have, try to call it on the regular dataframe (so not the sam
2 #if it takes way too long then call the function on the sampled version
3 #make sure to store these new tables into different variables we can use later on
4 sp500_df = pd.read_csv("/work/snp_500.csv")
5 climateCommentsDF_grouped = group_comments_by_sector(
6 climateCommentsDF, sp500_df, keep_single_sector_only=False
7)[0]
8 climateCommentsDF_grouped.to_csv("climateCommentsDF_grouped.csv", index=False)
```

```
1 #TODO call the function here on all of the dataframes we have, try to call it on the regular dataframe (so not the sam
2 #if it takes way too long then call the function on the sampled version
3 #make sure to store these new tables into different variables we can use later on
4 sp500_df = pd.read_csv("/work/snp_500.csv")
5 sustainabilityCommentsDF_grouped = group_comments_by_sector(
6 sustainabilityCommentsDF_sampled, sp500_df, keep_single_sector_only=False
7)[0]
8 sustainabilityCommentsDF_grouped.to_csv("sustainabilityCommentsDF_grouped.csv", index=False)
9 sustainabilityCommentsDF_grouped
```

# Six Step Sentiment Analysis Process (VADER)

...  
Install & Import Libraries,  
Load Reddit Data

Apply VADER to each  
reddit comment

Analyze Sentiment by  
Sector

Step 1

Step 2

Step 3

Step 4

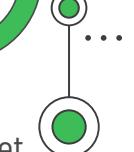
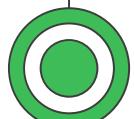
Step 5

Step 5

Get VADER sentiment  
scores

Convert compound score  
to sentiment label  
(positive, negative,  
neutral)

Save Final Dataset



...

```
1 !pip install nltk pandas
2
3 import pandas as pd
4 import nltk
5 from nltk.sentiment import SentimentIntensityAnalyzer
6
7 # Download the VADER lexicon (run once)
8 nltk.download("vader_lexicon")
9
```

# 1. Install and Import Libraries + Load Reddit Data



## 2. Get VADER sentiment scores

```
1 # Function that returns all VADER scores
2 def apply_vader(text):
3 if pd.isna(text):
4 text = ""
5 text = str(text)
6 scores = sia.polarity_scores(text)
7 return pd.Series({
8 "Negative Score": scores["neg"],
9 "Neutral Score": scores["neu"],
10 "Positive Score": scores["pos"],
11 "Compound Score": scores["compound"]
12 })
13
```



### 3. Apply VADER to each reddit comment

```
1 # Add vader sentiment analysis
2 climateCommentsDF_grouped[["Negative Score", "Neutral Score", "Positive Score", "Compound Score"]] = (
3 climateCommentsDF_grouped["Text"].apply(apply_vader)
4)
```

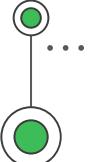


### 4. Convert compound score to sentiment label (positive, negative, neutral)

```
1 # convert compound score → discrete sentiment
2 def label_sentiment(compound):
3 if compound >= 0.05:
4 return "positive"
5 elif compound <= -0.05:
6 return "negative"
7 else:
8 return "neutral"
9
10 climateCommentsDF_grouped["vader_sentiment"] = climateCommentsDF_grouped["Compound Score"].apply(label_sentiment)
```



...



...



...

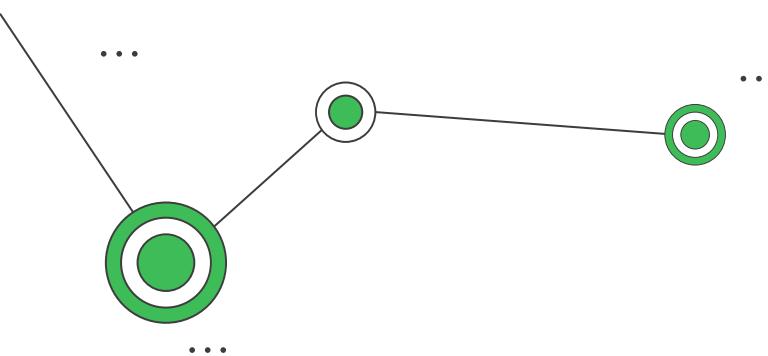
## 5. Analyze Sentiment by Sector

```
1 sector_sentiment_climate = (
2 climateCommentsDF_grouped
3 .groupby("mentioned_sector")["Compound Score"]
4 .mean()
5 .sort_values()
6)
7
8 sector_sentiment_climate
```

mentioned\_sector

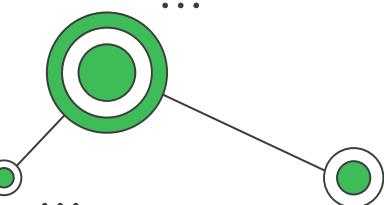
|                        |          |
|------------------------|----------|
| Real Estate            | 0.029970 |
| Industrials            | 0.137623 |
| Consumer Discretionary | 0.158428 |
| Consumer Staples       | 0.183451 |
| Communication Services | 0.211088 |
| Utilities              | 0.236732 |
| Materials              | 0.267280 |
| Technology             | 0.272308 |
| Financials             | 0.313255 |
| Energy                 | 0.328844 |
| Health Care            | 0.510546 |

Name: Compound Score, dtype: float64



## 6. Save Final Dataset!

```
climateCommentsDF_grouped.to_csv("climateCommentsDF_grouped_with_sentiment.csv", index=False)
```



# Repeat with other two subreddits...

Environment (envCommentsDF\_grouped):

```
1 envCommentsDF_grouped = pd.read_csv('/work/envCommentsDF_grouped.csv')
2 #envCommentsDF_grouped.head()
3 envCommentsDF_grouped[["Negative Score", "Neutral Score", "Positive Score", "Compound Score"]] = (
4 envCommentsDF_grouped["Text"].apply(apply_vader)
5)
6 envCommentsDF_grouped["vader_sentiment"] = envCommentsDF_grouped["Compound Score"].apply(label_sentiment)
7
8 sector_sentiment_env = (
9 envCommentsDF_grouped
10 .groupby("mentioned_sector")["Compound Score"]
11 .mean()
12 .sort_values()
13)
14
15 envCommentsDF_grouped.to_csv("envCommentsDF_grouped_with_sentiment.csv", index=False)
```

## Sustainability (sustainabilityCommentsDF\_grouped)

```
1 sustainabilityCommentsDF_grouped = pd.read_csv('/work/sustainabilityCommentsDF_grouped.csv')
2 sustainabilityCommentsDF_grouped[["Negative Score", "Neutral Score", "Positive Score", "Compound Score"]] = (
3 sustainabilityCommentsDF_grouped[["Text"]].apply(apply_vader)
4)
5 sustainabilityCommentsDF_grouped["vader_sentiment"] = sustainabilityCommentsDF_grouped["Compound Score"].apply(label_s
6
7 sector_sentiment_sustainability = (
8 sustainabilityCommentsDF_grouped
9 .groupby("mentioned_sector")["Compound Score"]
10 .mean()
11 .sort_values()
12)
13 sector_sentiment_sustainability
14 sustainabilityCommentsDF_grouped.to_csv("sustainabilityCommentsDF_grouped_with_sentiment.csv", index=False)
```

## 1 sector\_sentiment\_env

| mentioned_sector       |          |
|------------------------|----------|
| Industrials            | 0.113069 |
| Consumer Discretionary | 0.169247 |
| Consumer Staples       | 0.195384 |
| Real Estate            | 0.202945 |
| Communication Services | 0.232625 |
| Health Care            | 0.251750 |
| Utilities              | 0.282159 |
| Energy                 | 0.292724 |
| Materials              | 0.306743 |
| Technology             | 0.391694 |
| Financials             | 0.413507 |

## 1 sector\_sentiment\_sustainability

| mentioned_sector       | Score    |
|------------------------|----------|
| Financials             | 0.296598 |
| Technology             | 0.316769 |
| Consumer Staples       | 0.334345 |
| Energy                 | 0.381421 |
| Materials              | 0.386758 |
| Utilities              | 0.394625 |
| Consumer Discretionary | 0.412874 |
| Industrials            | 0.498292 |
| Communication Services | 0.655223 |
| Health Care            | 0.708948 |
| Real Estate            | 0.771938 |

# Results

## Cross-dataset Analysis

Most of the sentiments are positive.  
Sustainability has the highest positivity.

# Challenges

## Data Cleaning

Company-level matching is PROBLEMATIC:

1. Company names
2. Running time is fairly long

# Insight 1. In the Climate dataset, sentiment is positive overall, but varies sharply across sectors

```
1 sector_sentiment_climate = (
2 climateCommentsDF_grouped
3 .groupby("mentioned_sector")["Compound Score"]
4 .mean()
5 .sort_values()
6)
7
8 sector_sentiment_climate
✓
mentioned_sector
Real Estate 0.029970
Industrials 0.137623
Consumer Discretionary 0.158428
Consumer Staples 0.183451
Communication Services 0.211088
Utilities 0.236732
Materials 0.267280
Technology 0.272308
Financials 0.313255
Energy 0.328844
Health Care 0.510546
Name: Compound Score, dtype: float64
```

## Insight 2. In the Environment dataset, Technology and Financials remain positive, while Industrials are viewed less favorably

```
1 sector_sentiment_env
✓
mentioned_sector
Industrials 0.113069
Consumer Discretionary 0.169247
Consumer Staples 0.195384
Real Estate 0.202945
Communication Services 0.232625
Health Care 0.251750
Utilities 0.282159
Energy 0.292724
Materials 0.306743
Technology 0.391694
Financials 0.413507
Name: Compound Score, dtype: float64
```

## Insight 3. In the Sustainability dataset, sentiment becomes overwhelmingly positive

```
1 sector_sentiment_sustainability
mentioned_sector
Financials 0.296590
Technology 0.316769
Consumer Staples 0.334345
Energy 0.381421
Materials 0.386750
Utilities 0.394625
Consumer Discretionary 0.412874
Industrials 0.498292
Communication Services 0.655223
Health Care 0.700940
Real Estate 0.711938
Name: Compound Score, dtype: float64
```



## Insight 4. Recurring Trends & Interpretation

- a. Technology, Materials, and Utilities consistently score positively across climate, environment, and sustainability topics.  
→ progress and innovation.
  
- b. Industrials and Real Estate repeatedly appear on the lower end for climate/environment conversations.  
→ responsible for emissions or environmental risk.

# Thanks for Listening!

