

# Final Project Report: Sentiment Analysis on Movie Reviews

## Abstract

This project focuses on sentiment analysis of movie reviews using the Kaggle dataset. The dataset is derived from Pang and Lee's original movie review corpus and expanded upon by Socher et al. It contains sentiment-labeled phrases that are manually annotated into five categories: negative, somewhat negative, neutral, somewhat positive, and positive. The primary goal of this project is to build a sentiment classifier by leveraging various features and classification algorithms.

Comprehensive preprocessing, feature extraction, and classification experiments were conducted to identify the best-performing configurations.

## Dataset

The Kaggle movie review dataset was utilized for this project. The dataset consists of phrases from movie reviews labeled with sentiment scores ranging from 0 (most negative) to 4 (most positive). The training data 'train.tsv' was used, and it contains 156,060 phrases. A subset of 2,000 (limit set =2000) phrases was selected randomly for the experiments.

The sentiment labels include:

- 0: Negative
- 1: Somewhat Negative
- 2: Neutral
- 3: Somewhat Positive
- 4: Positive

## Part 1 : Data Preprocessing And Filtering→

The text preprocessing and filtering steps included the following:

- Converting text to lowercase.
- Removing punctuation using regular expressions.
- Tokenizing the text into individual words using NLTK.
- Removing stop-words using both default and custom lists tailored for the movie review context.
- Filtering tokens with fewer than three characters.

```
#Defining preprocess_text function
def preprocess_text(line):
    # Convert text to lowercase
    tokens = re.split(r'\s+', line.lower())
    # Remove punctuation
    punctuation_pattern = re.compile(r'![#$%&()*+,-./:;=>?@[\]^`{|}~]')
    cleaned_tokens = [punctuation_pattern.sub("", i) for i in tokens]
    # Combine default, custom, and additional stopwords
    stop_words = set(stopwords.words('english')).union(set(custom_stopwords)).union(set(additional_stopwords))
    # Remove stopwords
    filtered_tokens = [token for token in cleaned_tokens if token not in stop_words]
    processed_line = " ".join(filtered_tokens)
    return processed_line
```

The impact of preprocessing was visualized using histograms, word clouds, and frequency distributions.

### 1. Converting text to lowercase →

The texts were converted to lowercase.

```
# Convert text to lowercase
tokens = re.split(r'\s+', line.lower())
```

The above code processes a line of text by converting it to lowercase and splitting it into a list of words based on whitespace.

### 2. Removing punctuation using regular expressions →

Punctuation in each token is identified and removed by replacing the punctuation characters with an empty string during preprocessing. This ensures that the text is cleaned for analysis by eliminating unnecessary symbols.

```
# Remove punctuation
punctuation_pattern = re.compile(r'![#$%&()*+,-./:;=>?@[\]^`{|}~]')
cleaned_tokens = [punctuation_pattern.sub("", i) for i in tokens]
```

The above code removes punctuation from a list of tokens by using a regular expression to match punctuation characters and replacing them with an empty string. It creates a new list of cleaned tokens, ensuring the text is normalized for further processing.

### 3. Removing stop-words →

The NLTK default stop-words list has been expanded by incorporating two additional sets of stop-words: custom stop-words (e.g., filler words, contractions, and phrases) and domain-specific stop-words relevant to movie reviews (e.g., “movie,” “character,” “review”). This ensures more comprehensive filtering of irrelevant words during text preprocessing.

- Example additions:

**Custom Stop-words:** Words like “actually”, “usually”, “however.”

**Domain-Specific Stop-words:** Terms like “movie”, “character”, “cinema”.

```
# Import the default list of stopwords from NLTK
default_stopwords = nltk.corpus.stopwords.words('english')

# Custom stopwords
custom_stopwords = [
    "could", "would", "might", "must", "need", "shall", "wo", "sha", "y", "'s", "'d", "'ll",
    "'t", "'m", "'re", "'ve", "'n't", "'i", "not", "no", "cannot", "can't", "didn't", "doesn't", "isn't",
    "aren't", "wasn't", "weren't", "won't", "don't", "haven't", "hasn't", "hadn't", "shouldn't", "wouldn't",
    "actually", "also", "always", "even", "ever", "just", "really", "still", "yet", "however", "though",
    "nevertheless", "furthermore", "therefore", "otherwise", "meanwhile", "although", "thus", "hence", "indeed",
    "perhaps", "especially", "specifically", "usually", "often", "sometimes", "certainly", "typically", "mostly",
    "generally", "about", "above", "across", "after", "against", "among", "around", "before", "behind", "below",
    "beneath", "beside", "between", "beyond", "inside", "outside", "through", "under", "upon", "within", "without",
    "again", "almost", "already", "alone", "anyway", "anywhere", "elsewhere", "henceforth", "hereby", "herein",
    "hereafter", "hereupon", "thereby", "therein", "thereafter", "thereupon", "whereupon", "whereby", "whenever",
    "whatever", "whoever", "whosoever", "whomever", "ourselves", "yourselves", "itself", "themselves", "someone",
    "something", "anyone", "anything", "everyone", "everything", "nothing", "nobody", "somewhere", "everywhere", "lrb", "rrb"
]

# Add new stopwords for sentiment and movie reviews context
additional_stopwords = [
    "movie", "film", "scene", "story", "character", "plot", "performance", "review", "rating",
    "audience", "actor", "director", "cinema", "entertainment", "dialogue", "viewer", "show", "series"
]
```

### 4. Tokenization →

Each phrase in the dataset is tokenized into individual words using NLTK’s word\_tokenize function. This ensures consistent splitting based on whitespace, punctuation, and other delimiters.

```
tokens = nltk.word_tokenize(punctuation_pattern[0])
withoutpreprocessing.append((tokens, int(punctuation_pattern[1])))

punctuation_pattern[0] = preprocess_text(punctuation_pattern[0])
tokens = nltk.word_tokenize(punctuation_pattern[0])
withpreprocessing.append((tokens, int(punctuation_pattern[1])))
```

### Tokenization Without Preprocessing :

- `nltk.word_tokenize(punctuation_pattern[0]):`

Tokenizes the raw phrase (`punctuation_pattern[0]`) into individual words. This splits the phrase based on spaces, punctuation, or other delimiters.
- `withoutpreprocessing.append((tokens, int(punctuation_pattern[1]))):`

The tokenized words (`tokens`) and their corresponding sentiment label (`punctuation_pattern[1]`) are stored in the list `withoutpreprocessing` as a tuple (`(tokens, label)`) for later processing.

### Tokenization With Preprocessing:

- `punctuation_pattern[0] = preprocess_text(punctuation_pattern[0]):`

The phrase is passed through the `preprocess_text` function, which performs text cleaning such as: lowercasing, removing punctuation, removing stop-words.
- `nltk.word_tokenize(punctuation_pattern[0]):` Tokenizes the preprocessed phrase into individual words.
- `withpreprocessing.append((tokens, int(punctuation_pattern[1]))):`

The tokenized words from the preprocessed text (`tokens`) and their corresponding sentiment label (`punctuation_pattern[1]`) are stored in the list `withpreprocessing`.

## 5. Filtering word tokens →

A distinct function called `filter_tokens_by_length()` was developed to eliminate tokens from the list that had a length of **2 characters or less**. This was necessary because certain short words, such as “ll” and “wo,” were identified, which were deemed irrelevant in the context of sentiment analysis.

```
def filter_tokens_by_length(t):
    tokens=[]
    for n in t[0]:
        if len(n)>2:
            tokens.append(n)
    return (tokens,t[1])
```

This filtering step was applied after tokenization to ensure that the processed token list contained only meaningful and relevant words, thereby reducing noise and improving the overall quality of the feature set.

## Part 2 : Feature Engineering →

The goal is to transform the Kaggle movie review phrases into a structured format that can be used for classification. Feature engineering starts with basic bag-of-words features and expands to advanced features such as bigrams, trigrams, sentiment lexicons, and part-of-speech tagging.

Both filtered and unfiltered tokens and features have been generated.

Various features were extracted from the preprocessed data:

- Bag-of-Words (BoW) Features
- Unigram Features
- Bigram Features: High-scoring bigrams based on chi-square statistics.
- Trigram Features
- POS Features: Counts of nouns, verbs, adjectives, and adverbs.
- Sentiment Lexicon Features: Weak and strong positive/negative counts derived from subjectivity lexicons.
- LIWC Features: Counts of positive and negative words using the LIWC lexicon.
- Combined Features: Integration of unigrams, bigrams, POS counts, and sentiment lexicon features.

Bag-of-Words (BoW) Features :

- The `get_top_words` function identifies the most frequent words in the dataset.

```
def get_top_words(tokens,i):
    tokens = nltk.FreqDist(tokens)
    wf = [cleaned_tokens for (cleaned_tokens,c) in tokens.most_common(i)]
    return wf
```

Unigram Features :

- The `create_unigram_features` function generates a binary feature for each word in the top words:

```
def create_unigram_features(d,wf):
    df= set(d)
    f = {}
    for word in wf:
        f['V_%s' % word] = (word in df)
    return f
```

Bigram Features :

- Bigrams are two consecutive words

- Bigrams are extracted using NLTK's BigramCollocationFinder:

```
def extract_top_bigrams(wordlist,n):
    bigram_measure = nltk.collocations.BigramAssocMeasures()
    finder = BigramCollocationFinder.from_words(wordlist)
    finder.apply_freq_filter(2)
    b_features = finder.nbest(bigram_measure.chi_sq,4000)
    return b_features[:n]
```

Trigram Features :

- Trigrams are three consecutive words
- Trigrams are extracted below:

```
def extract_and_visualize_top_trigrams(tokens, top_n=10):

    # Extract trigrams
    trigram_measures = TrigramAssocMeasures()
    finder = TrigramCollocationFinder.from_words(tokens)
    finder.apply_freq_filter(2) # Minimum frequency of trigrams
    trigram_freq = finder.ngram_fd.items()
    sorted_trigrams = sorted(trigram_freq, key=lambda x: x[1], reverse=True)[:top_n]
```

POS Features :

- POS tagging assigns grammatical roles to words (e.g., noun, verb, adjective). Certain POS categories, like adjectives ("amazing") and adverbs ("very"), are strong indicators of sentiment.
- The extract\_pos\_features function counts nouns, verbs, adjectives, and adverbs in a phrase:

```
def extract_pos_features(document, word_features):
    document_words = set(document)
    tagged_words = nltk.pos_tag(document)
    features = {}
    for word in word_features:
        features['contains({})'.format(word)] = (word in document_words)
    numNoun = 0
    numVerb = 0
    numAdj = 0
    numAdverb = 0
    for (word, tag) in tagged_words:
        if tag.startswith('N'): numNoun += 1
        if tag.startswith('V'): numVerb += 1
        if tag.startswith('J'): numAdj += 1
        if tag.startswith('R'): numAdverb += 1
    features['nouns'] = numNoun
    features['verbs'] = numVerb
    features['adjectives'] = numAdj
    features['adverbs'] = numAdverb
    return features
```

Sentiment Lexicon Features :

- Sentiment lexicons are dictionaries of words labeled as positive, negative, or neutral. These features count the occurrences of such words in the text.
- Lexicons (e.g., Subjectivity Lexicon) are loaded, and positive/negative word counts are calculated:

```
def extract_subjectivity_features(document, word_features, SL):
    document_words = set(document)
    features = {}
    for word in word_features:
        features['V_{}'.format(word)] = (word in document_words)
    # count variables for the 4 classes of subjectivity
    weakPos = 0
    strongPos = 0
    weakNeg = 0
    strongNeg = 0
    for word in document_words:
        if word in SL:
            strength, posTag, isStemmed, polarity = SL[word]
            if strength == 'weaksubj' and polarity == 'positive':
                weakPos += 1
            if strength == 'strongsubj' and polarity == 'positive':
                strongPos += 1
            if strength == 'weaksubj' and polarity == 'negative':
                weakNeg += 1
            if strength == 'strongsubj' and polarity == 'negative':
                strongNeg += 1
        features['positivecount'] = weakPos + (2 * strongPos)
        features['negativecount'] = weakNeg + (2 * strongNeg)

    if 'positivecount' not in features:
        features['positivecount'] = 0
    if 'negativecount' not in features:
        features['negativecount'] = 0

    return features
```

LIWC Features :

- LIWC (Linguistic Inquiry and Word Count) is a lexicon-based feature extraction method that uses predefined dictionaries of words associated with positive and negative sentiment.
- Below is the code where I extracted LIWC features

```
def extract_liwc_features(doc, word_features, poslist, neglist):
    doc_words = set(doc)
    features = {}

    for word in word_features:
        features['contains({})'.format(word)] = (word in doc_words)

    pos = 0
    neg = 0
    for word in doc_words:
        if sentiment_read_LIWC_pos_neg_words.isPresent(word, poslist):
            pos += 1
        elif sentiment_read_LIWC_pos_neg_words.isPresent(word, neglist):
            neg += 1
        features['positivecount'] = pos
        features['negativecount'] = neg

    if 'positivecount' not in features:
        features['positivecount'] = 0
    if 'negativecount' not in features:
        features['negativecount'] = 0

    return features
```

## Combined Features :

- Combining multiple feature types (e.g., unigrams, bigrams, lexicons, POS tagging) to create a rich representation of the text.
- The extract\_combined\_features function integrates word features, bigrams, and lexicon-based counts:

```
def extract_combined_features(doc,word_features,SL,poslist,neglist):
    doc_words = set(doc)
    features={}

    for word in word_features:
        features['contains({})'.format(word)] = (word in doc_words)

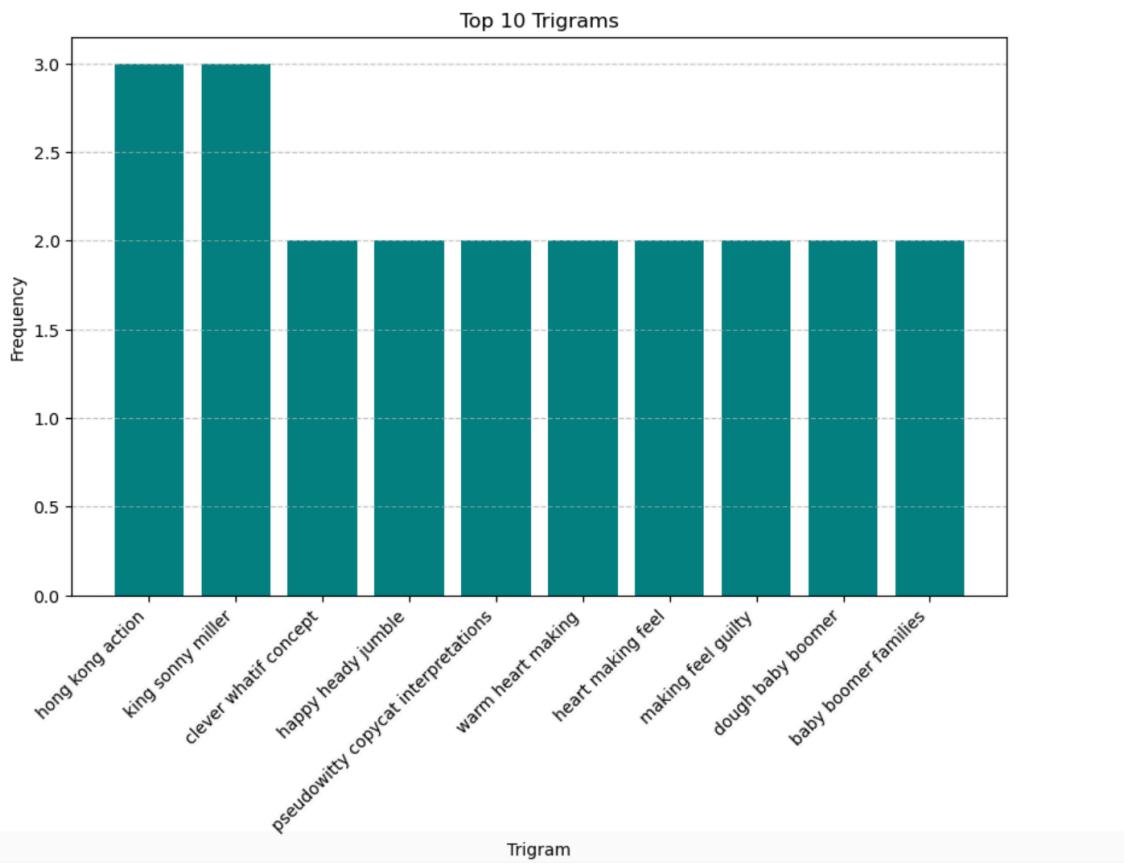
    weakPos = 0
    strongPos = 0
    weakNeg = 0
    strongNeg = 0
    for word in doc_words:
        if sentiment_read_LIWC_pos_neg_words.isPresent(word,poslist):
            strongPos +=1
        elif sentiment_read_LIWC_pos_neg_words.isPresent(word,neglist):
            strongNeg +=1
        elif word in SL:
            strength, posTag, isStemmed, polarity = SL[word]
            if strength == 'weaksubj' and polarity == 'positive':
                weakPos += 1
            if strength == 'strongsubj' and polarity == 'positive':
                strongPos += 1
            if strength == 'weaksubj' and polarity == 'negative':
                weakNeg += 1
            if strength == 'strongsubj' and polarity == 'negative':
                strongNeg += 1
        features['positivecount'] = weakPos + (2 * strongPos)
        features['negativecount'] = weakNeg + (2 * strongNeg)

    if 'positivecount' not in features:
        features['positivecount'] = 0
    if 'negativecount' not in features:
        features['negativecount'] = 0

    return features
```

## Part 3 : Data Visualization →

### 1. Trigrams :

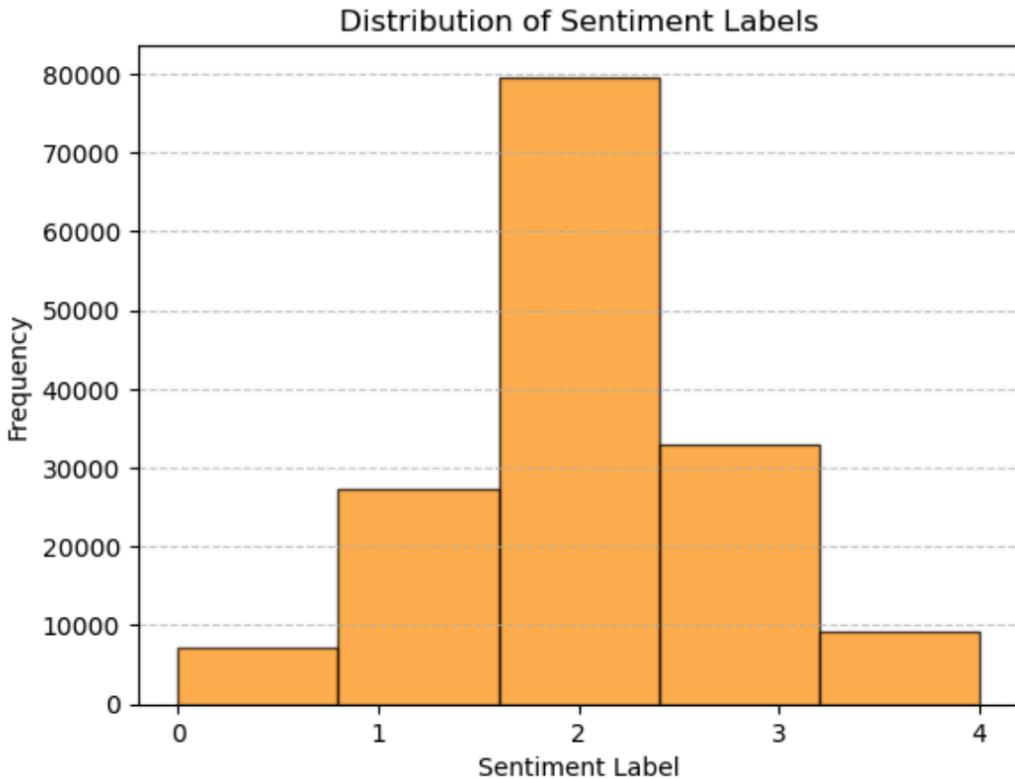


```
Top Trigrams with Frequencies: [(['hong', 'kong', 'action'), 3), (['king', 'sonny', 'miller'), 3), ([['clever', 'whatif', 'concept'), 2), (['happy', 'heady', 'jumble'), 2), ([['pseudowitty', 'copycat', 'interpretations'), 2), (['warm', 'heart', 'making'), 2), ([['heart', 'mak', 2), (['making', 'feel', 'guilty'), 2), (['dough', 'baby', 'boomer'), 2), ([['baby', 'boomer', 'families'), 2)]
```

1. **Top Trigrams:** The most frequent trigrams (e.g., “hong kong action,” “king sonny miller”) highlight recurring themes in the dataset, such as action genres and specific character references.
2. **Varied Frequency:** The trigrams show varied frequencies, with some appearing more often (e.g., 3 times), indicating their prominence in the dataset.
3. **Thematic Relevance:** Trigrams like “warm heart making” and “clever whatif concept” indicate emotional storytelling, unique plot dynamics, or creative descriptions.

4. **Cultural References:** Phrases such as “baby boomer families” reflect societal and demographic contexts within the reviews.

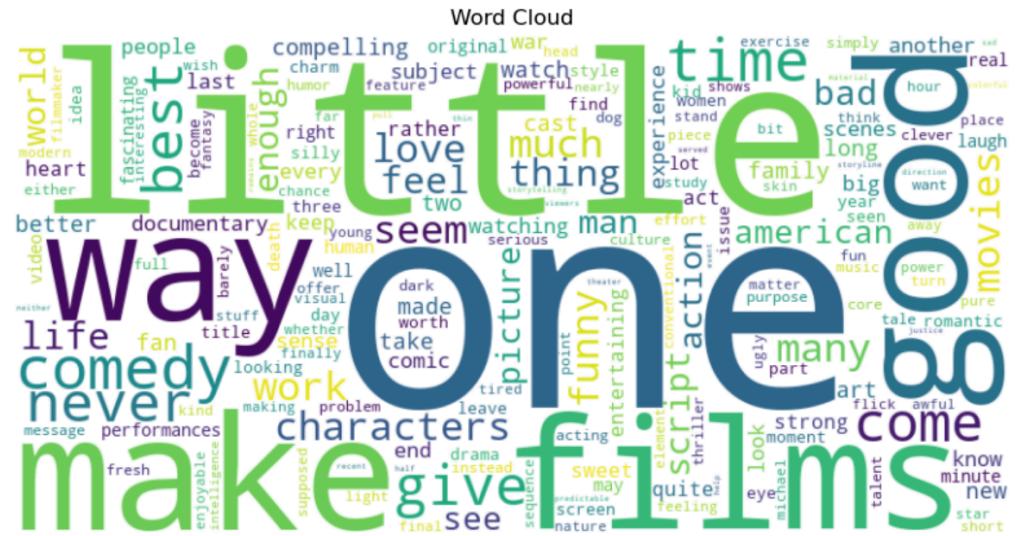
## 2. Histogram : Distribution of Sentiment Labels



The dataset contains a diverse range of sentiment labels, offering a comprehensive representation of movie reviews.

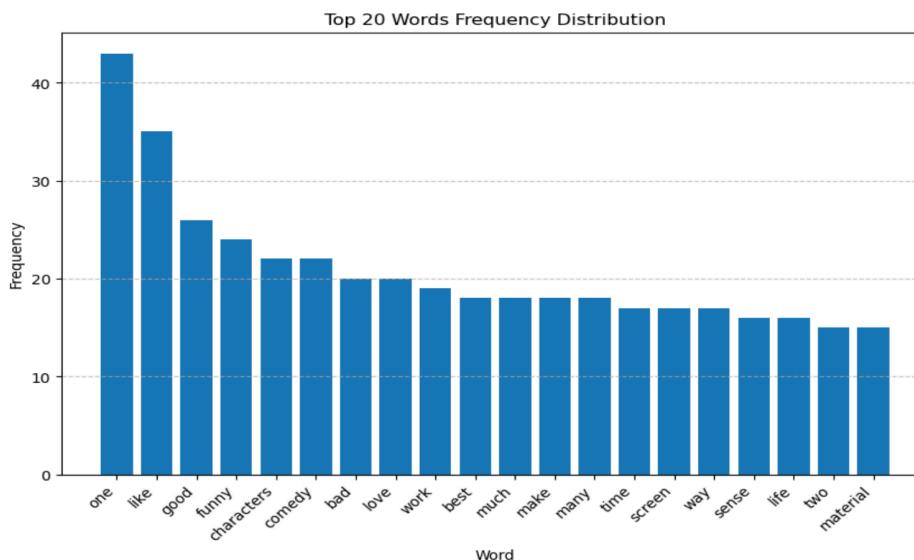
- Neutral sentiment (label 2) is well-represented, providing a strong foundation for understanding balanced opinions.
- Labels 1 (Negative) and 3 (Positive) have comparable frequencies, ensuring sufficient data for training these sentiment categories.

### 3. Word Cloud :



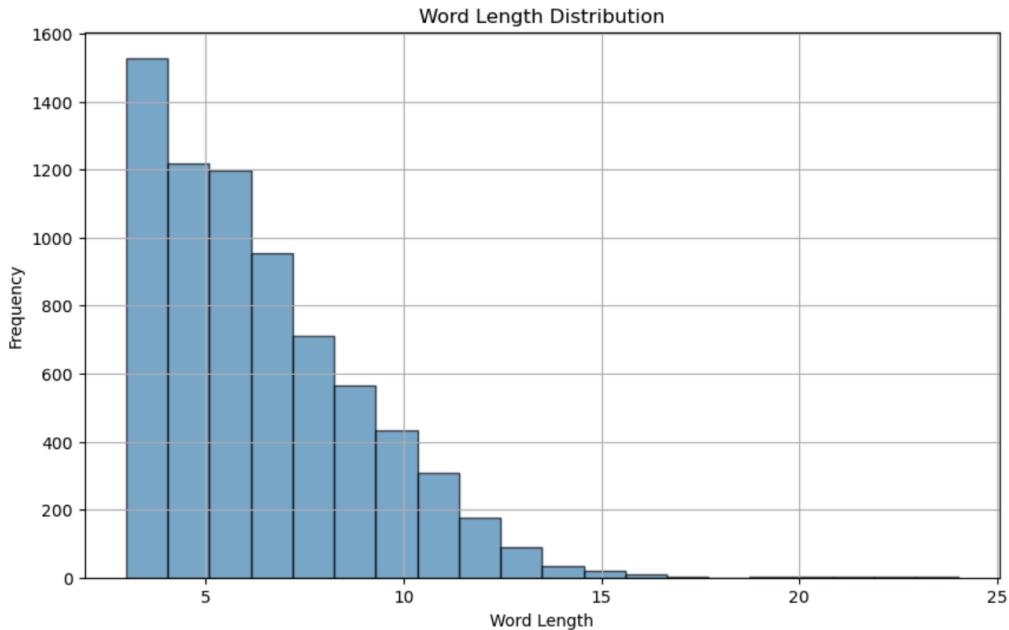
- The word cloud highlights frequent terms like “one,” “make,” “films,” and “way,” indicating their importance in the dataset.
- Words such as “comedy,” “characters,” and “funny” reflect recurring themes related to movie genres and storytelling.
- The presence of diverse terms shows the dataset’s focus on emotional, thematic, and cinematic aspects of reviews.

### 4. Histogram : Top 20 Words Frequency Distribution



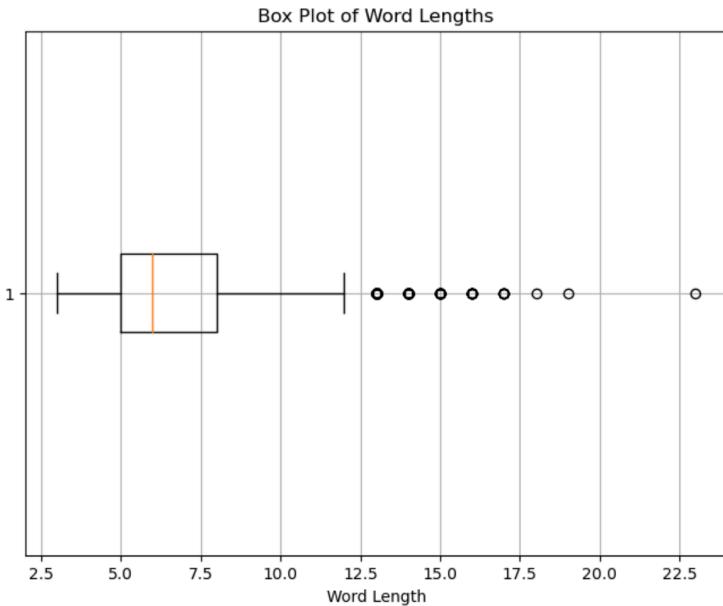
- The word “one” has the highest frequency, indicating its significant usage in movie reviews, possibly as a reference to singular experiences or descriptions.
- Words like “good,” “funny,” and “characters” are frequently used, reflecting common themes of humor and positive descriptions in the dataset.
- Sentiment-neutral words like “like” and “make” suggest structural or transitional language is prevalent in reviews.
- The distribution highlights a balance between descriptive terms (“best,” “bad,” “comedy”) and thematic references (“life,” “material”).

## 5. Histogram : Word Length Distribution



- Most words in the dataset have lengths between 4 and 7 characters, indicating the prevalence of medium-length descriptive terms.
- The frequency decreases sharply for longer words, suggesting fewer complex or multi-syllabic words are used in the reviews.
- This distribution aligns with typical natural language usage, where shorter, commonly used words dominate.

6. Boxplot :



1. **Median Word Length:** The orange line shows the median word length, approximately **5-6 characters**, indicating that most words in the dataset are of moderate length.
2. **Interquartile Range (IQR):** The box demonstrates that the majority of word lengths fall between **4 and 7 characters**, which is typical for natural language usage.
3. **Outliers:** The dots outside the whiskers represent words exceeding **15 characters**, likely rare or complex terms, proper nouns, or compound words.
4. **Dataset Variability:** The plot reveals that word lengths are fairly consistent, with minimal extreme variations, ensuring a balanced representation of short and medium-length words in the dataset.

**NOTE :**

- Both the box plot and histogram show that the majority of word lengths fall between **4 and 7 characters**, with a median around **5-6 characters**.
- The histogram's low-frequency bins for longer words align with the box plot's outliers, representing rare, unusually long words exceeding **15 characters**.
- Both visualizations confirm that the dataset predominantly contains medium-length words with minimal extreme variations, providing consistent insights into word length distribution.

## Part 4 : Experiments →

### 1. Experiment: Baseline with Unigram Features

- Used **unigram word features** (individual words) as the baseline for text classification.
- Extracted the most frequent words from the corpus and represented each phrase as a binary vector, indicating the presence or absence of these words.
- Function used: `create_unigram_features`.
- **Purpose:**  
To establish a baseline accuracy using simple bag-of-words features.
- **Results:**  
Unigram features provided reasonable accuracy but lacked contextual information.

### 2. Experiment: Stopword Filtering

- Preprocessed the text by removing common stopwords (e.g., “the,” “and”) to reduce noise.
  - Additionally, included a custom stopword list to remove context-specific irrelevant words (e.g., “movie,” “film”).
- Function used: `preprocess_text`.
- **Purpose:**  
To determine if removing irrelevant words improves classification accuracy.
  - **Results:**
    - Stop-word filtering improved model performance slightly by reducing vocabulary size and focusing on meaningful words.

### 3. Experiment: Sentiment Lexicon Features

- Used the **LIWC sentiment lexicon** to extract counts of positive and negative words in each phrase.
- Function used: `extract_liwc_features`.
- **Purpose:**  
To incorporate domain knowledge from pre-defined sentiment dictionaries.
- **Results:**  
LIWC features enhanced the model’s ability to identify strongly positive or negative sentiments, especially for extreme sentiment classes.

### 4. Experiment: Part-of-Speech (POS) Features

- Extracted counts of specific POS tags (e.g., nouns, verbs, adjectives, adverbs) for each phrase.
- Function used: `extract_pos_features`.
- **Purpose:**  
To capture the grammatical structure and sentiment-indicating words (e.g., adjectives like “great” or “terrible”).

- **Results:**  
POS features added context and improved performance for nuanced sentiments.

### 5. Experiment: Combined Features

- Combined multiple feature sets, including:
  - Unigrams
  - Bigrams
  - LIWC sentiment lexicon counts
  - POS features
  - Function used: extract\_combined\_features.
- **Purpose:**  
To create a richer representation of the text by combining multiple linguistic and semantic patterns.
- **Results:**  
Combined features provided the best classification performance, highlighting the importance of integrating different types of information.

### 6. Experiment: Varying Vocabulary Sizes

- Tested models with different sizes of the vocabulary (e.g., top 100, 500, 1000 words) to evaluate how the number of features impacts accuracy.
- Function used: get\_top\_words.
- **Purpose:**
- To optimize the trade-off between feature richness and computational efficiency.
- **Results:**
- Moderate vocabulary sizes (e.g., top 500 words) balanced performance and efficiency, while very large vocabularies introduced noise.

## Part 5 : Advance Experiments

### 1. Experiment: Using Different Classifiers

- Evaluated multiple classifiers with the feature sets, including:  
**Naïve Bayes:** NLTK's Naïve Bayes Classifier.  
**SVM:** Implemented using Sci-Kit Learn.  
**Random Forest:** Also implemented using Sci-Kit Learn.
- Functions used: `evaluate_naive_bayes_accuracy`,  
`evaluate_svm_classifier`, `evaluate_random_forest`.
- **Purpose:**  
To compare the effectiveness of different machine learning algorithms for sentiment classification.
- **Results:**  
SVM and Random Forest classifiers outperformed Naïve Bayes for richer feature sets, likely due to their ability to handle non-linear relationships.

### 2. Experiment: Cross-Validation

- Used **k-fold cross-validation** ( $k=5$ ) to evaluate the model's performance on different splits of the dataset.
- Function used: `cross_validation_PRF`.
- **Purpose:**  
To ensure robust evaluation by testing the model on multiple subsets of the data.
- **Results:**  
Cross-validation provided stable accuracy and highlighted consistent performance improvements with richer features.

```
def cross_validation_PRF(num_folds, featuresets, labels):  
    subset_size = int(len(featuresets)/num_folds)  
    print('Each fold size:', subset_size)  
    # for the number of labels - start the totals lists with zeroes  
    num_labels = len(labels)  
    total_precision_list = [0] * num_labels  
    total_recall_list = [0] * num_labels  
    total_F1_list = [0] * num_labels
```

### 3. Experiment: Lexicon Combination

- Combined LIWC features with another lexicon (e.g., Subjectivity Lexicon) to enhance sentiment detection.
- Function used: `extract_combined_features`.
- **Purpose:**  
To evaluate whether integrating multiple sentiment lexicons improves accuracy.
- **Results:**  
Combining lexicons provided a slight boost in accuracy, particularly for nuanced sentiments.

## Part 6 : Results →

The experiments compared various feature sets and classifiers. Key observations include:

- Unigrams and bigrams provided a strong baseline for classification.
- Combined features incorporating sentiment lexicons and POS tags significantly improved performance.
- The Random Forest classifier demonstrated the highest accuracy, followed by SVM.
- Preprocessed data consistently outperformed unfiltered data.

Explanation of these results is in the end.

### UNFILTERED :→

#### 1. Unigram Unfiltered :

Average Accuracy : 0.0985

	Average Precision	Recall	F1	Per Label
0	0.000	0.000	0.000	
1	0.175	0.291	0.217	
2	0.887	0.572	0.695	
3	0.185	0.357	0.244	
4	0.080	0.306	0.123	

	Macro Average Precision	Recall	F1	Over All Labels
	0.266	0.305	0.256	

Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}

	Micro Average Precision	Recall	F1	Over All Labels
	0.510	0.431	0.441	

Unigram Unfiltered :  
Each fold size: 400  
Fold 0

	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.159	0.302	0.208
2	0.921	0.570	0.704
3	0.190	0.364	0.250
4	0.034	0.250	0.061

Fold 1

	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.188	0.261	0.218
2	0.883	0.603	0.717
3	0.194	0.383	0.257
4	0.083	0.500	0.143

Fold 2

	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.203	0.271	0.232
2	0.861	0.582	0.695
3	0.168	0.356	0.229
4	0.038	0.111	0.057

Fold 3

	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.221	0.425	0.291
2	0.902	0.550	0.683
3	0.163	0.314	0.215
4	0.037	0.167	0.061

Fold 4

	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.107	0.195	0.138
2	0.869	0.553	0.676
3	0.212	0.367	0.269
4	0.208	0.500	0.294

Unigram Unfiltered :

Accuracy :  
0.515

	0	1	2	3	4	
0	<.>	4	3	1	.	
1	.	<4>34	5	.		
2	.	5<90>	2	.		
3	1	9	19	<8>	.	
4	.	2	4	8	<1>	

(row = reference; col = test)

2. Bigram Unfiltered :

Bigram Unfiltered :			
Each fold size: 400			
Fold 0			
	Precision	Recall	F1
0	0.333	0.172	0.227
1	0.134	0.324	0.190
2	0.847	0.622	0.717
3	0.250	0.339	0.288
4	0.207	0.375	0.267
Fold 1			
	Precision	Recall	F1
0	0.077	0.062	0.069
1	0.156	0.227	0.185
2	0.801	0.630	0.705
3	0.258	0.407	0.316
4	0.167	0.211	0.186
Fold 2			
	Precision	Recall	F1
0	0.143	0.118	0.129
1	0.188	0.261	0.218
2	0.796	0.613	0.693
3	0.168	0.296	0.215
4	0.154	0.182	0.167
Fold 3			
	Precision	Recall	F1
0	0.143	0.095	0.114
1	0.260	0.476	0.336
2	0.842	0.596	0.698
3	0.214	0.339	0.263
4	0.111	0.200	0.143
Fold 4			
	Precision	Recall	F1
0	0.160	0.250	0.195
1	0.133	0.200	0.160
2	0.791	0.594	0.679
3	0.282	0.375	0.322
4	0.167	0.250	0.200

Average Accuracy : 0.0965

	Average Precision	Recall	F1	Per Label
0	0.171	0.140	0.147	
1	0.174	0.298	0.218	
2	0.815	0.611	0.698	
3	0.235	0.351	0.281	
4	0.161	0.243	0.192	
Macro Average	Precision	Recall	F1	Over All Labels
	0.311	0.329	0.307	
Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}				
Micro Average	Precision	Recall	F1	Over All Labels
	0.499	0.452	0.461	

### Bigram Unfiltered :

Accuracy :

0.53

	0	1	2	3	4	
0	<2>	2	3	1	.	
1	1	<8>23	9	2		
2	1	7<86>	2	1		
3	6	8	15	<7>1		
4	2	1	1	8	<3>	

(row = reference; col = test)

### 3. POS Unfiltered :

Pos Unfiltered :			
Each fold size: 400			
Fold 0			
	Precision	Recall	F1
0	0.400	0.176	0.245
1	0.122	0.233	0.160
2	0.853	0.623	0.720
3	0.202	0.347	0.256
4	0.172	0.357	0.233
Fold 1			
	Precision	Recall	F1
0	0.308	0.174	0.222
1	0.188	0.250	0.214
2	0.811	0.635	0.712
3	0.183	0.347	0.239
4	0.167	0.235	0.195
Fold 2			
	Precision	Recall	F1
0	0.143	0.111	0.125
1	0.188	0.226	0.205
2	0.761	0.605	0.674
3	0.137	0.271	0.182
4	0.154	0.143	0.148
Fold 3			
	Precision	Recall	F1
0	0.143	0.087	0.108
1	0.156	0.279	0.200
2	0.837	0.583	0.688
3	0.163	0.291	0.209
4	0.111	0.200	0.143
Fold 4			
	Precision	Recall	F1
0	0.160	0.182	0.170
1	0.147	0.234	0.180
2	0.791	0.581	0.670
3	0.165	0.298	0.212
4	0.125	0.125	0.125

Average Accuracy : 0.0915

Average Precision	Recall	F1	Per Label
0	0.231	0.146	0.174
1	0.160	0.244	0.192
2	0.810	0.605	0.693
3	0.170	0.311	0.220
4	0.146	0.212	0.169

Macro Average Precision	Recall	F1	Over All Labels
0.303	0.304	0.289	

Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}			
Micro Average Precision	Recall	F1	
0.480	0.429	0.439	Over All Labels

Pos Unfiltered :

Accuracy :

0.51

	0	1	2	3	4	
0	<3>	2	3	.	.	
1	2	<7>25	7	2		
2	2	8<86>	1	.		
3	8	9	15	<5>	.	
4	3	1	3	7	<1>	

(row = reference; col = test)

4. SL Unfiltered :

SL Unfiltered :

Accuracy :

0.52

	0	1	2	3	4	
0	<3>	1	3	1	.	
1	1	<6>25	9	2		
2	2	7<85>	2	1		
3	5	8	16	<6>	2	
4	2	1	1	7	<4>	

(row = reference; col = test)

SL Unfiltered :			
Each fold size: 400			
Fold 0			
	Precision	Recall	F1
0	0.400	0.250	0.308
1	0.134	0.306	0.186
2	0.858	0.627	0.724
3	0.238	0.345	0.282
4	0.241	0.318	0.275
Fold 1			
	Precision	Recall	F1
0	0.154	0.125	0.138
1	0.188	0.240	0.211
2	0.796	0.638	0.708
3	0.280	0.448	0.344
4	0.167	0.211	0.186
Fold 2			
	Precision	Recall	F1
0	0.143	0.118	0.129
1	0.203	0.277	0.234
2	0.811	0.615	0.700
3	0.189	0.360	0.248
4	0.154	0.190	0.170
Fold 3			
	Precision	Recall	F1
0	0.143	0.118	0.129
1	0.234	0.419	0.300
2	0.837	0.586	0.689
3	0.184	0.295	0.226
4	0.111	0.188	0.140
Fold 4			
	Precision	Recall	F1
0	0.200	0.294	0.238
1	0.133	0.227	0.168
2	0.791	0.590	0.676
3	0.235	0.312	0.268
4	0.167	0.211	0.186

Average Accuracy : 0.095

Average Precision	Recall	F1	Per Label
0	0.208	0.181	0.188
1	0.178	0.294	0.220
2	0.818	0.611	0.699
3	0.225	0.352	0.274
4	0.168	0.223	0.191
Macro Average Precision	Recall	F1	Over All Labels
0.320	0.332	0.315	
Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}			
Micro Average Precision	Recall	F1	Over All Labels
0.501	0.452	0.462	

## 5. LIWC Unfiltered :

```

extract_liwc_features Unfiltered :
Each fold size: 400
Fold 0
    Precision      Recall      F1
0       0.333      0.167      0.222
1       0.122      0.263      0.167
2       0.842      0.620      0.714
3       0.226      0.328      0.268
4       0.207      0.375      0.267
Fold 1
    Precision      Recall      F1
0       0.308      0.182      0.229
1       0.156      0.227      0.185
2       0.816      0.639      0.716
3       0.258      0.414      0.318
4       0.125      0.231      0.162
Fold 2
    Precision      Recall      F1
0       0.143      0.133      0.138
1       0.203      0.283      0.236
2       0.796      0.606      0.688
3       0.189      0.353      0.247
4       0.192      0.208      0.200
Fold 3
    Precision      Recall      F1
0       0.143      0.111      0.125
1       0.247      0.463      0.322
2       0.842      0.581      0.687
3       0.184      0.333      0.237
4       0.185      0.250      0.213
Fold 4
    Precision      Recall      F1
0       0.120      0.188      0.146
1       0.133      0.192      0.157
2       0.796      0.596      0.682
3       0.247      0.350      0.290
4       0.125      0.176      0.146

```

Average Accuracy : 0.0945

	Average Precision	Recall	F1	Per Label
0	0.209	0.156	0.172	
1	0.172	0.286	0.214	
2	0.818	0.608	0.698	
3	0.221	0.356	0.272	
4	0.167	0.248	0.198	

	Macro Average Precision	Recall	F1	Over All Labels
	0.318	0.331	0.310	

Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}

	Micro Average Precision	Recall	F1	Over All Labels
	0.498	0.451	0.459	

```
extract_liwc_features Unfiltered :
```

```
Accuracy :  
0.535
```

	0	1	2	3	4	
0	<2>	3	2	1	.	
1	2	<8>24	7	2		
2	1	8<86>	2	.		
3	7	7	15	<8>	.	
4	2	1	1	8	<3>	

(row = reference; col = test)

## 6. Combined SL-LIWC Unfiltered :

```
Combined SL extract_liwc_features Unfiltered :  
Each fold size: 400
```

Fold 0			
	Precision	Recall	F1
0	0.333	0.172	0.227
1	0.134	0.324	0.190
2	0.847	0.622	0.717
3	0.250	0.339	0.288
4	0.207	0.375	0.267
Fold 1			
	Precision	Recall	F1
0	0.077	0.062	0.069
1	0.156	0.227	0.185
2	0.801	0.630	0.705
3	0.258	0.407	0.316
4	0.167	0.211	0.186
Fold 2			
	Precision	Recall	F1
0	0.143	0.118	0.129
1	0.188	0.261	0.218
2	0.796	0.613	0.693
3	0.168	0.296	0.215
4	0.154	0.182	0.167
Fold 3			
	Precision	Recall	F1
0	0.143	0.095	0.114
1	0.260	0.476	0.336
2	0.842	0.596	0.698
3	0.214	0.339	0.263
4	0.111	0.200	0.143
Fold 4			
	Precision	Recall	F1
0	0.160	0.250	0.195
1	0.133	0.200	0.160
2	0.791	0.594	0.679
3	0.282	0.375	0.322
4	0.167	0.250	0.200

Average Accuracy : 0.0965

	Average Precision	Recall	F1	Per Label
0	0.171	0.140	0.147	
1	0.174	0.298	0.218	
2	0.815	0.611	0.698	
3	0.235	0.351	0.281	
4	0.161	0.243	0.192	

	Macro Average Precision	Recall	F1	Over All Labels
	0.311	0.329	0.307	

Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}  
Micro Average Precision Recall F1 Over All Labels  
0.499 0.452 0.461

Combined SL extract\_liwc\_features Unfiltered :

Accuracy :

0.53

	0	1	2	3	4	
-----+-----						
0   <2>	2	3	1	.		
1   1	<8>	23	9	2		
2   1	7<86>	2	1			
3   6	8	15	<7>	1		
4   2	1	1	8	<3>		
-----+-----						

(row = reference; col = test)

FILTERED :→

1. Unigram :

Unigram filtered :			
Each fold size: 400			
Fold 0			
	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.073	0.250	0.113
2	0.926	0.537	0.680
3	0.179	0.341	0.234
4	0.034	0.333	0.062
Fold 1			
	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.125	0.320	0.180
2	0.917	0.568	0.701
3	0.172	0.410	0.242
4	0.000	0.000	0.000
Fold 2			
	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.062	0.154	0.089
2	0.910	0.553	0.688
3	0.158	0.417	0.229
4	0.000	0.000	0.000
Fold 3			
	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.156	0.462	0.233
2	0.951	0.530	0.681
3	0.122	0.308	0.175
4	0.037	0.200	0.062
Fold 4			
	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.093	0.233	0.133
2	0.901	0.520	0.659
3	0.141	0.387	0.207
4	0.167	0.500	0.250

Average Accuracy : 0.0975

Average Accuracy : 0.0975			
Average Precision	Recall	F1	Per Label
0	0.000	0.000	0.000
1	0.102	0.284	0.150
2	0.921	0.541	0.682
3	0.154	0.373	0.218
4	0.048	0.207	0.075
Macro Average Precision	Recall	F1	Over All Labels
0.245	0.281	0.225	
Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}			
Micro Average Precision	Recall	F1	Over All Labels
0.504	0.413	0.413	

Unigram filtered :

Accuracy :  
0.53

	0	1	2	3	4	
0	<.>	3	5	.	.	
1	.	<1>37	5	.	.	
2	.	3<93>	1	.	.	
3	.	3	22<11>	1	.	
4	.	.	4	10	<1>	

(row = reference; col = test)

## 2. Bigram :

Bigram filtered :			
Each fold size: 400			
Fold 0			
	Precision	Recall	F1
0	0.067	0.125	0.087
1	0.183	0.385	0.248
2	0.863	0.612	0.716
3	0.298	0.342	0.318
4	0.172	0.417	0.244
Fold 1			
	Precision	Recall	F1
0	0.077	0.125	0.095
1	0.094	0.146	0.114
2	0.777	0.615	0.687
3	0.301	0.373	0.333
4	0.083	0.125	0.100
Fold 2			
	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.188	0.293	0.229
2	0.821	0.587	0.685
3	0.179	0.304	0.225
4	0.115	0.200	0.146
Fold 3			
	Precision	Recall	F1
0	0.071	0.071	0.071
1	0.156	0.324	0.211
2	0.832	0.565	0.673
3	0.214	0.333	0.261
4	0.000	0.000	0.000
Fold 4			
	Precision	Recall	F1
0	0.160	0.667	0.258
1	0.120	0.250	0.162
2	0.780	0.569	0.658
3	0.318	0.342	0.329
4	0.083	0.118	0.098

Average Accuracy : 0.0955

	Average Precision	Recall	F1	Per Label
0	0.075	0.198	0.102	
1	0.148	0.280	0.193	
2	0.814	0.590	0.684	
3	0.262	0.339	0.293	
4	0.091	0.172	0.118	

	Macro Average Precision	Recall	F1	Over All Labels
	0.278	0.316	0.278	

	Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}		
		Micro Average Precision	Recall F1 Over All Labels
		0.491	0.433 0.446

Bigram filtered :

Accuracy :

0.49

	0	1	2	3	4	
0	<1>	3	2	2	.	
1	1	<5>28	9	.	.	
2	.	8<80>	9	.	.	
3	4	5	17	<9>2	.	
4	2	.	3	7	<3>	

(row = reference; col = test)

### 3. POS :

Pos filtered : Each fold size: 400			
Fold 0			
	Precision	Recall	F1
0	0.133	0.133	0.133
1	0.122	0.256	0.165
2	0.858	0.617	0.718
3	0.202	0.250	0.224
4	0.138	0.286	0.186
Fold 1			
	Precision	Recall	F1
0	0.077	0.091	0.083
1	0.125	0.216	0.158
2	0.811	0.612	0.697
3	0.280	0.394	0.327
4	0.042	0.077	0.054
Fold 2			
	Precision	Recall	F1
0	0.071	0.091	0.080
1	0.188	0.255	0.216
2	0.796	0.588	0.677
3	0.189	0.327	0.240
4	0.115	0.200	0.146
Fold 3			
	Precision	Recall	F1
0	0.071	0.067	0.069
1	0.195	0.395	0.261
2	0.864	0.568	0.685
3	0.112	0.224	0.150
4	0.000	0.000	0.000
Fold 4			
	Precision	Recall	F1
0	0.160	0.222	0.186
1	0.120	0.281	0.168
2	0.817	0.571	0.672
3	0.224	0.302	0.257
4	0.125	0.214	0.158

Average Accuracy : 0.0955

Average Precision	Recall	F1	Per Label
0	0.103	0.121	0.110
1	0.150	0.281	0.194
2	0.829	0.591	0.690
3	0.201	0.299	0.239
4	0.084	0.155	0.109

Macro Average Precision	Recall	F1	Over All Labels
0.273	0.290	0.268	

Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}

Micro Average Precision	Recall	F1	Over All Labels
0.486	0.421	0.436	

Pos filtered :

Accuracy :  
0.48

	0	1	2	3	4	
0	<1>	3	2	2	.	
1	2	<1>31	9	.	.	
2	1	5<86>	5	.	.	
3	5	7	18	<5>2	.	
4	2	1	3	6	<3>	

(row = reference; col = test)

#### 4. SL :

SL filtered :			
Each fold size: 400			
Fold 0			
	Precision	Recall	F1
0	0.200	0.500	0.286
1	0.220	0.450	0.295
2	0.858	0.617	0.718
3	0.298	0.321	0.309
4	0.172	0.417	0.244
Fold 1			
	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.125	0.174	0.145
2	0.796	0.631	0.704
3	0.301	0.389	0.339
4	0.167	0.333	0.222
Fold 2			
	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.141	0.257	0.182
2	0.816	0.580	0.678
3	0.179	0.288	0.221
4	0.154	0.250	0.190
Fold 3			
	Precision	Recall	F1
0	0.071	0.067	0.069
1	0.130	0.312	0.183
2	0.837	0.556	0.668
3	0.194	0.302	0.236
4	0.074	0.154	0.100
Fold 4			
	Precision	Recall	F1
0	0.120	0.375	0.182
1	0.120	0.257	0.164
2	0.780	0.562	0.654
3	0.282	0.316	0.298
4	0.167	0.250	0.200

Average Accuracy : 0.0945

	Average Precision	Recall	F1	Per Label
0	0.078	0.188	0.107	
1	0.147	0.290	0.194	
2	0.817	0.589	0.684	
3	0.251	0.323	0.281	
4	0.147	0.281	0.191	

Macro Average	Precision	Recall	F1	Over All Labels
	0.288	0.334	0.291	

Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}

Micro Average	Precision	Recall	F1	Over All Labels
	0.494	0.438	0.448	

---

SL filtered :

Accuracy :  
0.49

	0	1	2	3	4	
0	<1>	3	2	2	.	
1	1	<6>27	9	.	.	
2	.	6<80>11	.	.	.	
3	3	5	19	<8>2	.	
4	2	.	3	7	<3>	

(row = reference; col = test)

5. LIWC :

```

extract_liwc_features filtered :
Each fold size: 400
Fold 0
    Precision      Recall      F1
0       0.200      0.250      0.222
1       0.122      0.263      0.167
2       0.842      0.588      0.693
3       0.298      0.352      0.323
4       0.103      0.429      0.167
Fold 1
    Precision      Recall      F1
0       0.077      0.100      0.087
1       0.125      0.200      0.154
2       0.835      0.625      0.715
3       0.290      0.458      0.355
4       0.083      0.125      0.100
Fold 2
    Precision      Recall      F1
0       0.000      0.000      0.000
1       0.156      0.278      0.200
2       0.806      0.574      0.671
3       0.189      0.277      0.225
4       0.077      0.167      0.105
Fold 3
    Precision      Recall      F1
0       0.000      0.000      0.000
1       0.130      0.333      0.187
2       0.848      0.547      0.665
3       0.194      0.311      0.239
4       0.037      0.077      0.050
Fold 4
    Precision      Recall      F1
0       0.160      0.500      0.242
1       0.160      0.333      0.216
2       0.806      0.568      0.667
3       0.271      0.319      0.293
4       0.125      0.231      0.162

```

Average Accuracy : 0.098

	Average Precision	Recall	F1	Per Label
0	0.087	0.170	0.110	
1	0.139	0.282	0.185	
2	0.827	0.581	0.682	
3	0.248	0.344	0.287	
4	0.085	0.206	0.117	

	Macro Average Precision	Recall	F1	Over All Labels
	0.277	0.316	0.276	

	Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}			
	Micro Average Precision	Recall	F1	Over All Labels
	0.493	0.432	0.442	

```
extract_liwc_features filtered :
```

```
Accuracy :
```

```
0.5
```

	0	1	2	3	4	
0	<3>	1	2	2	.	
1	1	<4>31	7	.	.	
2	.	8<81>	8	.	.	
3	3	5	17<10>	2	.	
4	2	.	4	7	<2>	

(row = reference; col = test)

## 6. Combined SI-LIWC:

```
Combined SL extract_liwc_features filtered :
```

```
Accuracy :
```

```
0.49
```

	0	1	2	3	4	
0	<1>	3	2	2	.	
1	1	<5>28	9	.	.	
2	.	8<80>	9	.	.	
3	4	5	17	<9>	2	
4	2	.	3	7	<3>	

(row = reference; col = test)

Combined SL extract\_liwc\_features filtered:  
 Each fold size: 400  
 Fold 0

	Precision	Recall	F1
0	0.067	0.125	0.087
1	0.183	0.385	0.248
2	0.863	0.612	0.716
3	0.298	0.342	0.318
4	0.172	0.417	0.244

Fold 1

	Precision	Recall	F1
0	0.077	0.125	0.095
1	0.094	0.146	0.114
2	0.777	0.615	0.687
3	0.301	0.373	0.333
4	0.083	0.125	0.100

Fold 2

	Precision	Recall	F1
0	0.000	0.000	0.000
1	0.188	0.293	0.229
2	0.821	0.587	0.685
3	0.179	0.304	0.225
4	0.115	0.200	0.146

Fold 3

	Precision	Recall	F1
0	0.071	0.071	0.071
1	0.156	0.324	0.211
2	0.832	0.565	0.673
3	0.214	0.333	0.261
4	0.000	0.000	0.000

Fold 4

	Precision	Recall	F1
0	0.160	0.667	0.258
1	0.120	0.250	0.162
2	0.780	0.569	0.658
3	0.318	0.342	0.329
4	0.083	0.118	0.098

Average Accuracy : 0.0955

	Average Precision	Recall	F1	Per Label
0	0.075	0.198	0.102	
1	0.148	0.280	0.193	
2	0.814	0.590	0.684	
3	0.262	0.339	0.293	
4	0.091	0.172	0.118	

	Macro Average Precision	Recall	F1	Over All Labels
	0.278	0.316	0.278	

Label Counts {0: 81, 1: 362, 2: 972, 3: 455, 4: 130}

	Micro Average Precision	Recall	F1	Over All Labels
	0.491	0.433	0.446	

Evaluation of Decision Tree :→

- **Unfiltered Features:**
  - **Best Accuracy:** 49.5% using **Subjectivity Lexicon (SL)** features.
  - **Worst Accuracy:** 41% using **POS features**, showing the limited impact of grammatical structure alone.
- **Filtered Features:**
  - **Best Accuracy:** 51% using **POS features**, indicating that filtering enhances the impact of POS features by removing noise.
  - **Worst Accuracy:** 46.5% using **Unigram features**, suggesting that unigram-based models are less effective after filtering.
- **Key Observations:**
  - Accuracy generally improves with filtering, as seen in POS features (41% to 51%).
  - Combining features like LIWC and SL doesn't always outperform individual features, likely due to overlapping information.

### 1. Unfiltered →

```
Unigram Unfiltered :  
Classifier-DecisionTree  
  
Accuracy : 0.46  
  
Bigram Unfiltered :  
Classifier-DecisionTree  
  
Accuracy : 0.435  
  
Pos Unfiltered :  
Classifier-DecisionTree  
  
Accuracy : 0.41  
  
SL Unfiltered :  
Classifier-DecisionTree  
  
Accuracy : 0.495  
  
extract_liwc_features Unfiltered :  
Classifier-DecisionTree  
  
Accuracy : 0.435  
  
Combined SL extract_liwc_features Unfiltered :  
Classifier-DecisionTree  
  
Accuracy : 0.45
```

## 2. Unfiltered →

```
Unigram filtered :  
Classifier-DecisionTree  
  
Accuracy : 0.465  
  
Bigram filtered :  
Classifier-DecisionTree  
  
Accuracy : 0.505  
  
Pos filtered :  
Classifier-DecisionTree  
  
Accuracy : 0.51  
  
SL filtered :  
Classifier-DecisionTree  
  
Accuracy : 0.49  
  
extract_liwc_features filtered :  
Classifier-DecisionTree  
  
Accuracy : 0.49  
  
Combined SL extract_liwc_features filtered :  
Classifier-DecisionTree  
  
Accuracy : 0.49
```

## Evaluation of SVM :

- **Unfiltered Features:**
  - **Best Accuracy:** 52% with **Bigram features** and **Combined LIWC+SL features**.
  - **Worst Accuracy:** 50% with **Unigram features**, showing a consistent baseline but limited improvement.
- **Filtered Features:**
  - **Best Accuracy:** 55% with **Unigram features** and **POS features**, indicating that filtering enhances the SVM's ability to leverage these features.
  - **Worst Accuracy:** 49.5% with **Filtered LIWC features**, suggesting LIWC's standalone impact diminishes after filtering.
- **Key Observations:**
  - SVM consistently outperforms Decision Trees in almost all cases.
  - Filtering significantly improves unigram and POS-based accuracy for SVM, highlighting its sensitivity to feature refinement.

## 1. UNFILTERED :

```
Unigram Unfiltered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.5  
  
Bigram Unfiltered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.52  
  
Pos Unfiltered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.505  
  
SL Unfiltered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.505  
  
extract_liwc_features Unfiltered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.515  
  
Combined SL extract_liwc_features Unfiltered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.52
```

## 2. FILTERED :

```
Unigram filtered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.55  
  
Bigram filtered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.51  
  
Pos filtered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.55  
  
SL filtered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.505  
  
extract_liwc_features filtered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.495  
  
Combined SL extract_liwc_features filtered :  
Classifier-evaluate_svm_classifier  
Accuracy : 0.51
```

Evaluation of Random Forest :

- **Unfiltered Features:**
  - **Best Accuracy:** 50% with **Unigram features**, showing Random Forest's baseline performance.
  - **Worst Accuracy:** 47% with **POS features** and **Bigram features**, indicating challenges in leveraging complex structures without preprocessing.
- **Filtered Features:**
  - **Best Accuracy:** 54% with **POS features**, showing the highest performance for Random Forest when features are filtered.
  - **Worst Accuracy:** 47.5% with **Filtered LIWC features**, indicating limited standalone utility of LIWC features for this classifier.
- **Key Observations:**
  - Random Forest performs best with filtered POS features, leveraging their ability to capture grammatical structure.
  - Filtered features consistently outperform unfiltered ones for Random Forest, reinforcing the importance of preprocessing.

## 1. FILTERED :

```
Unigram Unfiltered :  
Classifier - Random Forest  
  
Accuracy: 0.5  
  
Bigram Unfiltered :  
Classifier - Random Forest  
  
Accuracy: 0.47  
  
Pos Unfiltered :  
Classifier - Random Forest  
  
Accuracy: 0.47  
  
SL Unfiltered :  
Classifier - Random Forest  
  
Accuracy: 0.485  
  
extract_liwc_features Unfiltered :  
Classifier - Random Forest  
  
Accuracy: 0.475  
  
Combined SL extract_liwc_features Unfiltered :  
Classifier - Random Forest  
  
Accuracy: 0.47
```

## 2. FILTERED :

```
Unigram filtered :  
Classifier - Random Forest  
Accuracy: 0.505  
  
Bigram filtered :  
Classifier - Random Forest  
Accuracy: 0.505  
  
Pos filtered :  
Classifier - Random Forest  
Accuracy: 0.54  
  
SL filtered :  
Classifier - Random Forest  
Accuracy: 0.505  
  
extract_liwc_features filtered :  
Classifier - Random Forest  
Accuracy: 0.475  
  
Combined SL extract_liwc_features filtered :  
Classifier - Random Forest  
Accuracy: 0.515
```

## Comparison Across Classifiers →

- **Overall Best Accuracy:**
  - 55% achieved by **SVM with Unigram and POS features (filtered)**.
- **Performance Hierarchy:**
  - SVM > Random Forest > Decision Tree, with SVM consistently achieving higher accuracy across most feature sets.
- **Impact of Filtering:**
  - Filtering consistently improves accuracy across all classifiers and feature sets.
  - For example, Decision Tree accuracy with POS features increases from 41% (unfiltered) to 51% (filtered).

## Feature-Specific Insights →

- **Unigrams:**
  - Performs consistently well across all classifiers and improves significantly with filtering, as seen in SVM (50% → 55%).
- **Bigrams:**
  - Moderate performance across classifiers; best accuracy (52%) is achieved with unfiltered SVM.

- **POS Features:**
  - Highly impactful for all classifiers, especially with filtering (e.g., 51% for Decision Tree, 55% for SVM, and 54% for Random Forest).
- **LIWC and Combined Features:**
  - While LIWC features contribute to improved accuracy, combining them with SL features doesn't always result in a significant boost, likely due to feature redundancy.

## Part 7 : Observations →

### Observations for the Report

1. **Best Classifier Performance:**
  - **SVM** performed the best overall, achieving the highest accuracy of **55%** with filtered Unigram and POS features, showcasing its strength in handling text data.
2. **Impact of Preprocessing (Filtering):**
  - Filtering consistently improved accuracy across all classifiers, especially for Decision Tree and Random Forest.
  - For example, Decision Tree accuracy with POS features increased from **41% (unfiltered)** to **51% (filtered)**.
3. **Effectiveness of POS Features:**
  - POS features outperformed other feature sets across classifiers, achieving **54% accuracy** with Random Forest and **55% accuracy** with SVM (both filtered).
4. **Unigram Features as a Strong Baseline:**
  - Unigrams showed reliable performance across classifiers, particularly after filtering, achieving **55% accuracy** with SVM.
5. **Limited Impact of Combined Features:**
  - Combining features like LIWC and SL did not significantly improve accuracy, likely due to feature redundancy.
6. **Classifiers Performance Hierarchy:**
  - SVM consistently outperformed Decision Tree and Random Forest across most feature sets, demonstrating its ability to handle nuanced patterns in text data.

**7. Bigram Features:**

- Bigrams showed moderate performance, with a maximum accuracy of **52%** (SVM, unfiltered), indicating limited added value compared to unigrams.

**8. LIWC Features:**

- LIWC features had marginal contributions to accuracy, with the highest accuracy of **51.5%** (Random Forest, filtered).

**9. Overall Performance:**

- Despite the improvements, the maximum accuracy of **55%** suggests potential for further feature enhancements, such as TF-IDF or embeddings.