Preprocessing: This step involves cleaning the text data by removing noise, such as punctuation, stopwords, and special characters, as well as performing tokenization, lemmatization, and stemming.

1. Feature Extraction:
   - TF-IDF (Term Frequency-Inverse Document Frequency): TF-IDF represents the importance of a word in a document relative to a collection of documents. It assigns weights to words based on their frequency in the document and their rarity across the entire corpus.
   - GloVe (Global Vectors for Word Representation): GloVe is a word embedding technique that represents words as dense vectors in a continuous vector space. It captures semantic relationships between words based on their co-occurrence statistics in a large corpus.
   - BERT (Bidirectional Encoder Representations from Transformers): BERT is a contextual word embedding technique that generates word embeddings by considering the surrounding context of each word. It captures complex semantic relationships and syntactic structures in the text data.

2. Sentiment Score Calculation:
   - VADER: Use the VADER sentiment analysis tool to calculate sentiment scores for each token or sentence in your text data. VADER provides positive, negative, and neutral sentiment scores for each token, along with an overall compound score that represents the aggregated sentiment.
   - SentiWordNet: Look up each token in the SentiWordNet lexicon to retrieve its sentiment scores (positive, negative, and objective scores). Calculate the sentiment score for each sentence or document by aggregating the sentiment scores of its constituent tokens.

3. Aggregation:
   - VADER: Optionally, aggregate the sentiment scores of individual tokens to obtain sentiment scores for entire sentences or documents. You can use simple averaging or weighted averaging based on token importance.
   - SentiWordNet: Similarly, aggregate the sentiment scores of individual tokens to obtain sentiment scores for sentences or documents.

4. Thresholding (Optional):
   - Depending on your requirements, you may choose to apply thresholding to the sentiment scores obtained from VADER or SentiWordNet to classify the sentiment as positive, negative, or neutral. For example, you could classify a document as "positive" if its compound sentiment score from VADER is above a certain threshold, and "negative" if it's below another threshold.

5. Analysis and Interpretation:
   - Once you have obtained sentiment scores for your text data, analyze and interpret the results. This may involve summarizing the distribution of sentiment scores, identifying patterns or trends in sentiment across different documents or topics, and understanding the overall sentiment of your dataset.

6. Visualization (Optional):
   - To aid in interpretation, you may choose to visualize the sentiment analysis results using plots or charts. For example, you could create histograms or pie charts to visualize the distribution of sentiment labels (positive, negative, neutral) in your dataset.

By following these steps, you can effectively perform sentiment analysis using lexicon-based models like VADER and SentiWordNet on your cleaned text data. Remember to validate the results and fine-tune any parameters or thresholds based on the specific requirements of your analysis.