## Review Length Class

```python
## Author : Goh Boon Xiang
from pyspark.sql.functions import size, col
import matplotlib.pyplot as plt

class ReviewLengthAnalyzer:
    def __init__(self, df_tokenized):
        """
        Initializes the ReviewLengthAnalyzer with the tokenized DataFrame.
        :param df_tokenized: Spark DataFrame containing tokenized reviews.
        """
        self.df_tokenized = df_tokenized

    def calculate_review_lengths(self):
        """
        Calculate the length of each review based on the number of tokens.
        :return: DataFrame with review lengths added.
        """
        df_review_length = self.df_tokenized.withColumn("review_length",
size(col("tokens")))
        return df_review_length

    def collect_review_lengths(self, df_review_length):
        """
        Collect the review lengths from the DataFrame as a list.
        :param df_review_length: DataFrame with review lengths.
        :return: List of review lengths.
        """
        return df_review_length.select("review_length").rdd.flatMap(lambda x:
x).collect()

    def plot_review_length_distribution(self, review_lengths):
        """
        Plot a histogram of review lengths.
        :param review_lengths: List of review lengths.
        """
        plt.figure(figsize=(10, 6))
        plt.hist(review_lengths, bins=30, color='mistyrose', edgecolor='black')
        plt.xlabel('Number of Tokens (Review Length)')
        plt.ylabel('Frequency')
        plt.title('Distribution of Review Lengths (Number of Tokens)')
        plt.show()
```