```
import nltk
```

```
nltk.download('stopwords')
```

```
    [nltk_data] Downloading package stopwords to /root/nltk_data...
    [nltk_data]   Unzipping corpora/stopwords.zip.
    True
```

```
nltk.download('punkt')
```

```
    [nltk_data] Downloading package punkt to /root/nltk_data...
    [nltk_data]   Unzipping tokenizers/punkt.zip.
    True
```

```
nltk.download('wordnet')
```

```
    [nltk_data] Downloading package wordnet to /root/nltk_data...
    True
```

```
import pandas as pd
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
import string

# Load the CSV file
df = pd.read_csv('/content/sample_data/Tag_dataset.csv')

# Data Cleaning
df.dropna(subset=['Description'], inplace=True)  # Drop rows with missing descriptions
df.drop_duplicates(subset=['Description'], keep='first', inplace=True)  # Remove duplicates

# Text Preprocessing
stop_words = set(stopwords.words('english'))
lemmatizer = WordNetLemmatizer()

def preprocess_text(text):
    tokens = word_tokenize(text.lower())  # Tokenization and lowercasing
    tokens = [word for word in tokens if word.isalnum()]  # Remove special characters
    tokens = [word for word in tokens if word not in stop_words]  # Remove stopwords
    tokens = [lemmatizer.lemmatize(word) for word in tokens]  # Lemmatization
    return tokens  # Return list of tokens instead of joined string

df['clean_description_tokens'] = df['Description'].apply(preprocess_text)  # Create new column with tokenized words

# Save the preprocessed data to a new CSV file
df.to_csv('preprocessed_data.csv', index=False)
```

```
nltk.download('averaged_perceptron_tagger')

    [nltk_data] Downloading package averaged_perceptron_tagger to
    [nltk_data]     /root/nltk_data...
    [nltk_data]   Unzipping taggers/averaged_perceptron_tagger.zip.
    True
```