





```
import pandas as pd

df = pd.read_pickle('consumer_complaint_dataset.data', compression='gzip')
```

df



| | topic | input | |
|--------|---|---|---|
| 0 | Debt collection | transworld systems inc. \nis trying to collect... |  |
| 1 | Credit reporting, credit repair services, or o... | I would like to request the suppression of the... |  |
| 2 | Debt collection | Over the past 2 weeks, I have been receiving e... |  |
| 3 | Credit reporting, credit repair services, or o... | I HAD FILED WITH CFPB ON XX/XX/XXXX19 TO HAVE ... | |
| 4 | Credit reporting, credit repair services, or o... | I have several accounts that the balance is in... | |
| ... | ... | ... | |
| 492250 | Consumer Loan | I was on automatic payment for my car loan. In... | |
| 492251 | Debt collection | I recieved a collections call from an unknown ... | |
| 492252 | Mortgage | On XXXX XXXX, 2015, I contacted XXXX XXXX, who... | |
| 492253 | Mortgage | I can not get from chase who services my mortg... | |
| 492254 | Credit card | I made a payment to CITI XXXX Credit Card on X... | |

492255 rows x 3 columns

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

# Download necessary NLTK resources
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')

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

def preprocess_text(text):
    # Convert to lowercase
    text = text.lower()
    # Tokenize the text
    words = word_tokenize(text)
    # Remove stopwords and punctuation, and lemmatize
    words = [lemmatizer.lemmatize(word) for word in words if word not in stop_words and word not in string.punctuation]
    return words

# Apply the pre-processing to the 'input' column
df['processed_input'] = df['input'].apply(preprocess_text)
```

```

[ nltk_data] Downloading package punkt to /root/nltk_data...
[ nltk_data] Package punkt is already up-to-date!
[ nltk_data] Downloading package stopwords to /root/nltk_data...
[ nltk_data] Package stopwords is already up-to-date!
[ nltk_data] Downloading package wordnet to /root/nltk_data...
[ nltk_data] Package wordnet is already up-to-date!

```

df

| | topic | input | processed_input |
|--------|---|---|---|
| 0 | Debt collection | transworld systems inc. \nis trying to collect... | [transworld, system, inc., trying, collect, de... |
| 1 | Credit reporting, credit repair services, or o... | I would like to request the suppression of the... | [would, like, request, suppression, following,... |
| 2 | Debt collection | Over the past 2 weeks, I have been receiving e... | [past, 2, week, receiving, excessive, amount, ... |
| 3 | Credit reporting, credit repair services, or o... | I HAD FILED WITH CFPB ON XX/XX/XXXX19 TO HAVE ... | [filed, cfpb, xx/xx/xxxx19, listed, account, d... |
| 4 | Credit reporting, credit repair services, or o... | I have several accounts that the balance is in... | [several, account, balance, incorrect, couple,... |
| ... | ... | ... | ... |
| 492250 | Consumer Loan | I was on automatic payment for my car loan. In... | [automatic, payment, car, loan, fine, print, s... |
| 492251 | Debt collection | I recieved a collections call from an unknown ... | [recieved, collection, call, unknown, company,... |
| 492252 | Mortgage | On XXXX XXXX, 2015, I contacted XXXX XXXX, who... | [xxxx, xxxx, 2015, contacted, xxxx, xxxx, bran... |
| 492253 | Mortgage | I can not get from chase who services my mortg... | [get, chase, service, mortgage, owns, original... |
| 492254 | Credit card | I made a payment to CITI XXXX Credit Card on X... | [made, payment, citi, xxxx, credit, card, xxxx... |

492255 rows × 3 columns

```

from gensim.models import Word2Vec

# Train Word2Vec on the processed text
model = Word2Vec(sentences=df['processed_input'], vector_size=100, window=5, min_count=2, workers=4)

# Save the model for later use
model.save("word2vec_complaints.model")

similar_words_debt = model.wv.most_similar('debt', topn=10)
similar_words_collection = model.wv.most_similar('collection', topn=10)
similar_words_risk = model.wv.most_similar('risk', topn=10)

print("Similar words to 'debt':", similar_words_debt)
print("Similar words to 'collection':", similar_words_collection)
print("Similar words to 'risk':", similar_words_risk)

```

⇒ Similar words to 'debt': [('debt-', 0.7821791768074036), ('deb', 0.7581514120101929), ('debt"', 0.71793
Similar words to 'collection': [('aargon', 0.7778722643852234), ('colection', 0.7114332914352417), ('th
Similar words to 'risk': [('danger', 0.64796382188797), ('exposure', 0.6130653619766235), ('reputational', 0.6130653619766235)]

```
from sklearn.manifold import TSNE
import matplotlib.pyplot as plt
import numpy as np # Import numpy

# Get the vectors of the similar words
words = ['debt', 'collection', 'risk'] + [word for word, _ in similar_words_debt + similar_words_collection + similar_words_risk]
word_vectors = [model.wv[word] for word in words]

# Convert word_vectors to a NumPy array
word_vectors = np.array(word_vectors) # Convert the list of vectors to a 2D numpy array

# Apply t-SNE to reduce to 2 dimensions
tsne = TSNE(n_components=2, random_state=42)
word_vec_tsne = tsne.fit_transform(word_vectors)

# Plotting the words
plt.figure(figsize=(10, 10))
plt.scatter(word_vec_tsne[:, 0], word_vec_tsne[:, 1])

for i, word in enumerate(words):
    plt.annotate(word, (word_vec_tsne[i, 0], word_vec_tsne[i, 1]))

plt.title("t-SNE Visualization of Word Embeddings")
plt.show()
```

