```
1 # Mount Google Drive to access dataset
2 from google.colab import drive
3 drive.mount('/content/drive')
```

```
   Mounted at /content/drive
```

```
1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 import nltk
```

```
1 nltk.download('stopwords')
2 nltk.download('punkt')
```

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

```
1 from nltk.corpus import stopwords
2 from nltk.tokenize import word_tokenize
3 from nltk.stem import PorterStemmer
4
5 from gensim.models import word2vec
6 from gensim.models import Word2Vec
7
8 from sklearn.feature_extraction.text import CountVectorizer
9 from sklearn.metrics.pairwise import cosine_similarity
```

```
 1 def pprocess(text):
 2     # Load the NLTK stop words
 3     stop_words = set(stopwords.words('english'))
 4     # Initialize the NLTK stemmer
 5     stemmer = PorterStemmer()
 6     tokens = word_tokenize(text)
 7
 8     # Lowercase the tokens
 9     tokens = [token.lower() for token in tokens]
10
11     # Remove punctuation and special characters
12     tokens = [token for token in tokens if token.isalpha()]
13
14     # Remove stop words
15     tokens = [token for token in tokens if token not in stop_words]
16
17     # Stem the tokens
18     tokens = [stemmer.stem(token) for token in tokens]
19
20     return tokens
```

```
 1 def preprocess(text):
 2     ret = []
 3     # Load the NLTK stop words
 4     stop_words = set(stopwords.words('english'))
 5     # Initialize the NLTK stemmer
 6     stemmer = PorterStemmer()
 7
 8     for tex in text:
 9         # Tokenize the sentence
10         tokens = word_tokenize(tex)
11
12         # Lowercase the tokens
13         tokens = [token.lower() for token in tokens]
14
15         # Remove punctuation and special characters
16         tokens = [token for token in tokens if token.isalpha()]
17
18         # Remove stop words
19         tokens = [token for token in tokens if token not in stop_words]
20
21         # Stem the tokens
22         tokens = [stemmer.stem(token) for token in tokens]
23         ret.append(tokens)
24
25     return ret
```

```
1 # Load the dataset into a DataFrame
2 df_0 = pd.read_csv('/content/drive/MyDrive/bick/train00.csv')
3 df_1 = pd.read_csv('/content/drive/MyDrive/bick/train01.csv')
4 df_2 = pd.read_csv('/content/drive/MyDrive/bick/train02.csv')
5 df_3 = pd.read_csv('/content/drive/MyDrive/bick/train03.csv')
```

```
 1 train_df0 = df_0['Review Text'].to_frame()
 2 train_df0.columns = ['text']
 3
 4 train_df1 = df_1['Review'].to_frame()
 5 train_df1.columns = ['text']
 6
 7 train_df2 = df_2['Review'].to_frame()
 8 train_df2.columns = ['text']
 9
10 train_df3 = df_3['review_text'].to_frame()
11 train_df3.columns = ['text']
12
13 train_df = pd.concat([train_df0, train_df1, train_df2, train_df3], ignore_index=True)
14 #train_df['corpus'] = ''
```

```
1 train_df
```

|       | text |
|-------|------|
| 0     | Absolutely wonderful - silky and sexy and comf… |
| 1     | Love this dress! it's sooo pretty. i happene… |
| 2     | I had such high hopes for this dress and reall… |
| 3     | I love, love, love this jumpsuit. it's fun, fl… |
| 4     | This shirt is very flattering to all due to th… |
| …     | … |
| 70453 | I oot this dress in the blue. it fits great--h… |
| 70454 | I was very patient with this dress. i was wait… |
| 70455 | The deep v doesn't gape, and flatters the neck… |
| 70456 | I saw this dress online this morning, went int… |
| 70457 | Super cute jacket .perfect for fall i can't st… |

70458 rows × 1 columns

```
1 train_df['text'].isna().sum()
```

```
2535
```

```
1 train_df = train_df.dropna(axis=0)
2 train_df = train_df.drop_duplicates(['text'], keep='first')
3 train_df = train_df.reset_index(drop=True)
```

```
1 train_df['corpus'] = train_df['text'].apply(pprocess)
```

```
1 train_df
```

```
1 corpus = train_df['corpus'].to_list()
```
        Absolutely wonderful silky and sexy and comf...        [absolut, wonder, silki, sexi, comfort]

```
1 # Train a Word2Vec model on the preprocessed corpus
2 model = Word2Vec(corpus, min_count=1, vector_size=100)
```

```
 1 # Compute the similarity between the input sentence and each sentence in the text file
 2 input_sentence = "As you put on the shirt, you suddenly realize that it feels like someone else's skin clinging uncomfortab
 3 input_tokens = pprocess(input_sentence)
 4 similar_sentences = []
 5 for i in range(len(corpus)):
 6     if len(corpus[i]) == 0 or len(input_tokens) == 0:
 7         continue
 8     similarity = model.wv.n_similarity(input_tokens, corpus[i])
 9     similar_sentences.append((train_df.loc[i, 'text'], similarity))
10
11 # Sort the similar sentences by similarity score in descending order and output the top n sentences
12 n = 30
13 top_similar_sentences = sorted(similar_sentences, key=lambda x: x[1], reverse=True)[:n]
14 tmp = 1
15 for sentence, similarity in top_similar_sentences:
16     print(f"\n{tmp}) sentence: \n\t", sentence, "\nsimilarity: ", similarity)
17     tmp += 1
```

```
   1) sentence:
          Lovely sweater. it fits baggy and form fitted somehow. i like the multi-fabric look. however, the main body
   similarity:  0.8000628

   2) sentence:
          I got this 25% off before christmas. i immediately started wearing it, thinking it would be a staple. i like
   similarity:  0.7970948

   3) sentence:
          Nice fabric, great color, would be lovely ... on someone with a straight-and-narrow body type. for ladies wi
   similarity:  0.7910863

   4) sentence:
          Nice fabric, magnificent color, would be lovely ... on someone with a straight-and-narrow body type. for lad
   similarity:  0.79052174

   5) sentence:
          I was hoping this clothe was a bit thicker, but when i tried it on the fabric was very thin and sheer. becau
   similarity:  0.7869623

   6) sentence:
          I was hoping this dress was a bit thicker, but when i tried it on the fabric was very thin and sheer. becaus
   similarity:  0.7840034

   7) sentence:
          I thought this was love at first sight. the color and design are both stunning. and it drapes just perfectly
   similarity:  0.7836629

   8) sentence:
          I loved the idea of this clothe, but i find that it isn't very flattering. i think the arm holes are actuall

   the material is soft, but a bit stuffy and i could tell it would give me crazy static in my hair.
   similarity:  0.7835511

   9) sentence:
          I loved the idea of this dress, but i find that it isn't very flattering. i think the arm holes are actually

   the material is soft, but a bit stuffy and i could tell it would give me crazy static in my hair.
   similarity:  0.78336215

  10) sentence:
          Wow. bailey 44 knocks it out of the park with this sensual, i-am-woman-hear-me-purr body-slimming clothe. we
   similarity:  0.77883756

  11) sentence:
          This top is really beautiful, but looks completely horrible on my 5'9" apple shaped body. in the words of my
   similarity:  0.77763677

  12) sentence:
          The material was thin, showed every lump you'd rather hide and clung to the body. even the sleeves were tigh
   similarity:  0.77623314

  13) sentence:
          I thought this was love at first sight. the color and design are both stunning. and it drapes just excellent
   similarity:  0.7758187
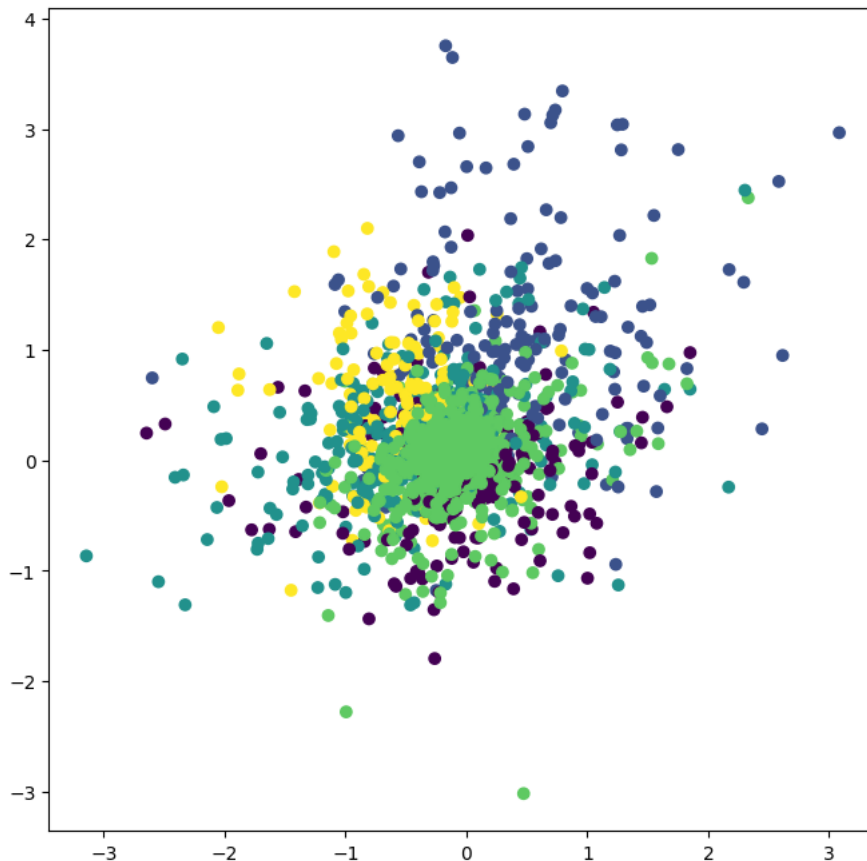```

```
1 import gensim
2 import numpy as np
```

```
3  import matplotlib.pyplot as plt
4  from sklearn.cluster import KMeans
5
6  # Load trained word2vec model
7  # model = gensim.models.Word2Vec.load('path/to/trained/model')
8
9  # Get embedding matrix and list of words
10 embeddings = model.wv.vectors
11 words = model.wv.index_to_key
12
13 # Cluster embeddings using K-Means
14 kmeans = KMeans(n_clusters=5, random_state=42)
15 clusters = kmeans.fit_predict(embeddings)
16
17 # Plot the clusters using the first two dimensions of the embeddings
18 plt.figure(figsize=(8, 8))
19 plt.scatter(embeddings[:, 0], embeddings[:, 1], c=clusters, cmap='viridis')
20 plt.show()
21
```
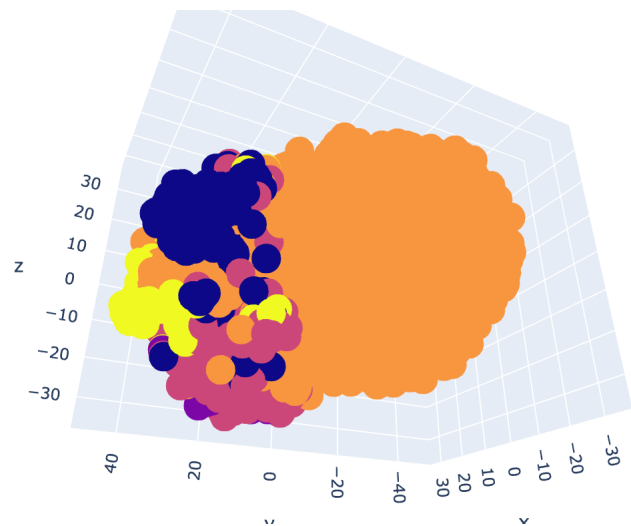
```
/usr/local/lib/python3.9/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will
  warnings.warn(
```



```
1  import gensim
2  import plotly.express as px
3  from sklearn.manifold import TSNE
4
5  # Load trained word2vec model
6  # model = gensim.models.Word2Vec.load('path/to/trained/model')
7
8  # Get embedding matrix and list of words
9  embeddings = model.wv.vectors
10 words = model.wv.index_to_key
11
12 # Use t-SNE to project embeddings into 2D space
13 tsne = TSNE(n_components=3, random_state=42)
14 embeddings_3d = tsne.fit_transform(embeddings)
15
16 # Create a dataframe with the 2D embeddings and their corresponding words
17 df = pd.DataFrame({'x': embeddings_3d[:, 0], 'y': embeddings_3d[:, 1], 'z': embeddings_3d[:, 2], 'word': words})
18
19 # Use Plotly to create a scatter plot of the embeddings
20 fig = px.scatter_3d(df, x='x', y='y', z='z', color=clusters, hover_data=['word'])
21 fig.show()
```
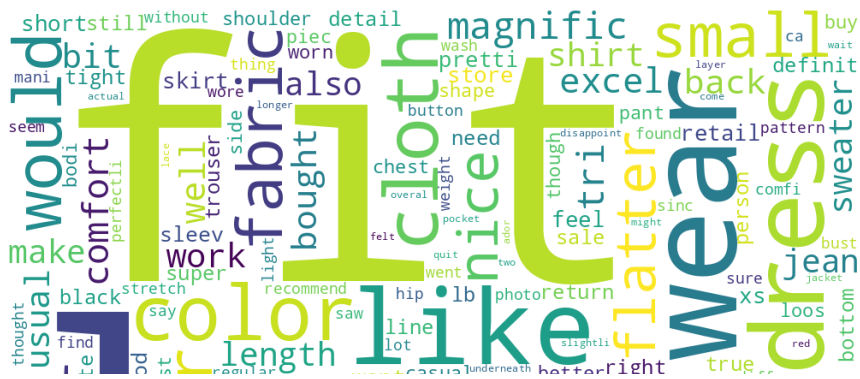
```
1  # Get the top 100 most common words
2  top_words = model.wv.index_to_key[:500]
3
4  # Assign frequency values based on how often each word appears in the model
5
6  word_frequencies = {}
7  for word in top_words:
8      word_frequencies[word] = model.wv.get_vecattr(word, "count")
9
10 from wordcloud import WordCloud
11
12 # Create the word cloud object
13 wordcloud = WordCloud(width=1000, height=800, background_color="white")
14
15 # Generate the word cloud
16 wordcloud.generate_from_frequencies(word_frequencies)
17
18 # Display the word cloud
19 import matplotlib.pyplot as plt
20 plt.figure(figsize=(10,8), facecolor=None)
21 plt.imshow(wordcloud, interpolation='bilinear')
22 plt.axis("off")
23 plt.tight_layout(pad=0)
24 plt.show()
25
```

1



✓　6초　　오후 6:58에 완료됨