

2048010

NLP - Case Study

E-Commerce Clothing Review Classification with TF

```
In [1]: 1 import os
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import seaborn as sns
5 import plotly.graph_objects as go
6 from wordcloud import WordCloud
7 from sklearn.model_selection import train_test_split
8 import numpy as np
9 import re
10 import plotly.express as px
11 from sklearn.metrics import confusion_matrix, accuracy_score
12 from sklearn.metrics import precision_score
13 from sklearn.metrics import recall_score
14 import nltk
15 from nltk.tokenize import word_tokenize
16 from nltk.corpus import stopwords
17 from nltk.stem import WordNetLemmatizer
18 from tensorflow.keras.preprocessing.text import Tokenizer
19 from tensorflow.keras.preprocessing.sequence import pad_sequences
20 from keras.layers import Dense, Dropout
21 from keras.layers import LSTM
22 from keras.models import Sequential
23 from keras.layers import Bidirectional
24 from keras.layers import Embedding
25 from keras.layers import GlobalAvgPool1D
26 import tensorflow as tf
27
28
29 for dirname, _, filenames in os.walk('/kaggle/input'):
30     for filename in filenames:
31         print(os.path.join(dirname, filename))
```

C:\ProgramData\Anaconda3\lib\site-packages\statsmodels\tools_testing.py:19: FutureWarning: pandas.util.testing is deprecated. Use the functions in the public API at pandas.testing instead.
import pandas.util.testing as tm

```
In [2]: 1 data = pd.read_csv("E:/Christ University/Semester-3/04_NLP/Case_Study/data/Womens Clothing E-Commerce R
2 data.head()
```

Out[2]:

	Clothing ID	Age	Title	Review Text	Rating	Recommended IND	Positive Feedback Count	Division Name	Department Name	Class Name
0	767	33	NaN	Absolutely wonderful - silky and sexy and comf...	4	1	0	Initmates	Intimate	Intimates
1	1080	34	NaN	Love this dress! it's sooo pretty. i happene...	5	1	4	General	Dresses	Dresses
2	1077	60	Some major design flaws	I had such high hopes for this dress and reall...	3	0	0	General	Dresses	Dresses
3	1049	50	My favorite buy!	I love, love, love this jumpsuit. it's fun, fl...	5	1	0	General Petite	Bottoms	Pants
4	847	47	Flattering shirt	This shirt is very flattering to all due to th...	5	1	6	General	Tops	Blouses

As you can see from the head of the dataset, we have some unnecessary features such as Clothing ID, Title. First of all, I will drop this features.

```
In [3]: 1 data = data.drop(['Title', 'Clothing ID', 'Positive Feedback Count'], axis=1)
        2 data.head()
```

Out[3]:

	Age	Review Text	Rating	Recommended IND	Division Name	Department Name	Class Name
0	33	Absolutely wonderful - silky and sexy and comf...	4	1	Intimates	Intimate	Intimates
1	34	Love this dress! it's sooo pretty. i happene...	5	1	General	Dresses	Dresses
2	60	I had such high hopes for this dress and reall...	3	0	General	Dresses	Dresses
3	50	I love, love, love this jumpsuit. it's fun, fl...	5	1	General Petite	Bottoms	Pants
4	47	This shirt is very flattering to all due to th...	5	1	General	Tops	Blouses

```
In [4]: 1 # Checking for the missing values
        2 count_NaN = data.isna().sum()
        3 count_NaN
```

```
Out[4]: Age          0
Review Text      845
Rating          0
Recommended IND  0
Division Name    14
Department Name  14
Class Name       14
dtype: int64
```

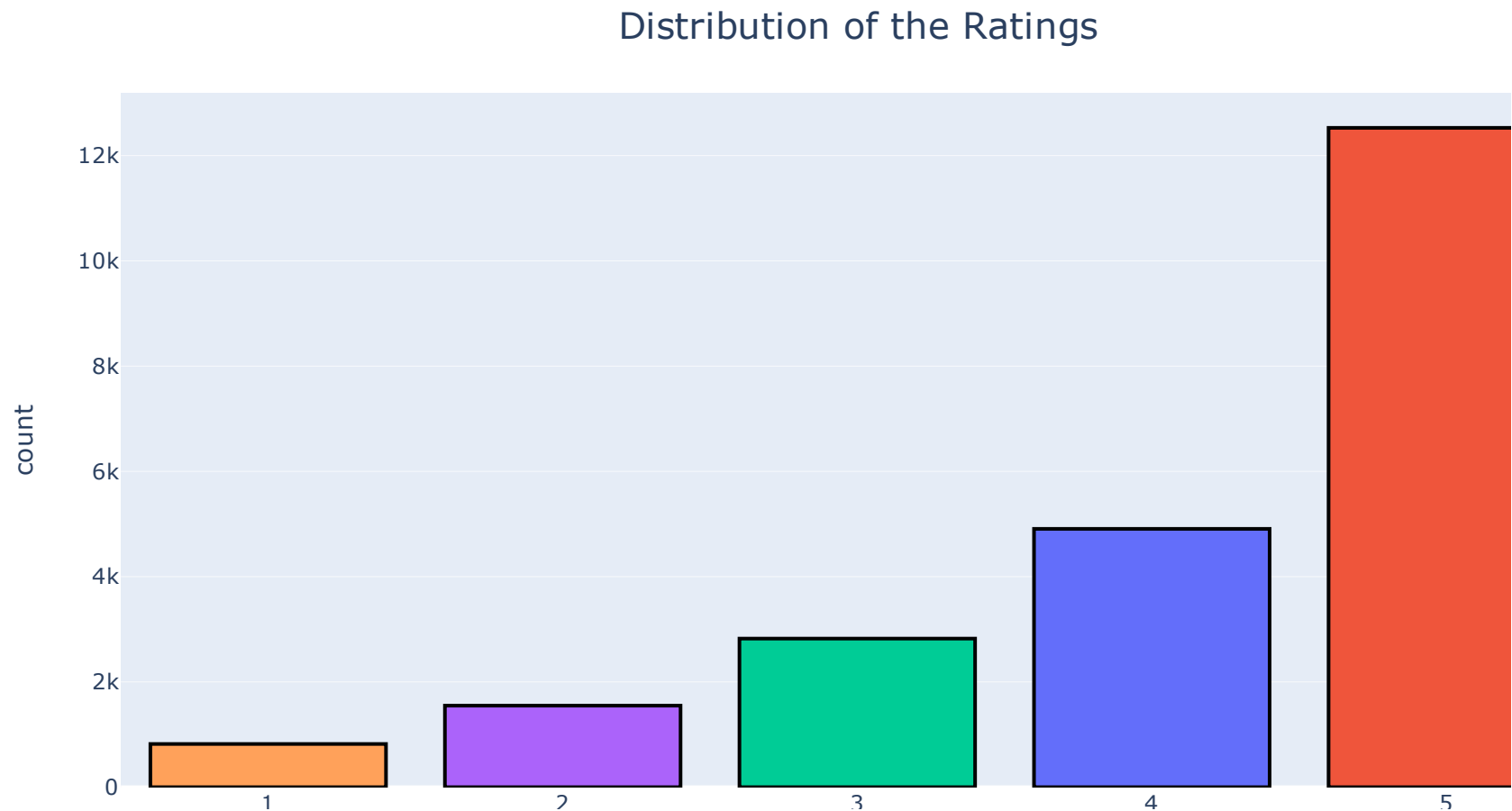
We will drop the missing rows from the dataset.

```
In [5]: 1 # Dropping the missing values in the rows
2 data = data.dropna(subset=['Review Text', 'Division Name', 'Department Name', 'Class Name'], axis=0)
3 data = data.reset_index(drop=True)
4
5 # Checking for the missing values after the drops
6 count_NaN_updated = data.isna().sum()
7 count_NaN_updated
```

```
Out[5]: Age      0
Review Text    0
Rating         0
Recommended IND 0
Division Name  0
Department Name 0
Class Name     0
dtype: int64
```

Distribution of the Ratings

```
In [6]: 1 fig = px.histogram(data['Rating'],
2                     labels={'value': 'Rating',
3                           'count': 'Frequency',
4                           'color': 'Rating'}, color=data['Rating'])
5 fig.update_layout(bargap=0.2)
6 fig.update_traces(marker=dict(line=dict(color='#000000', width=2)))
7 fig.update_layout(title_text='Distribution of the Ratings',
8                   title_x=0.5, title_font=dict(size=20))
9 fig.show()
```

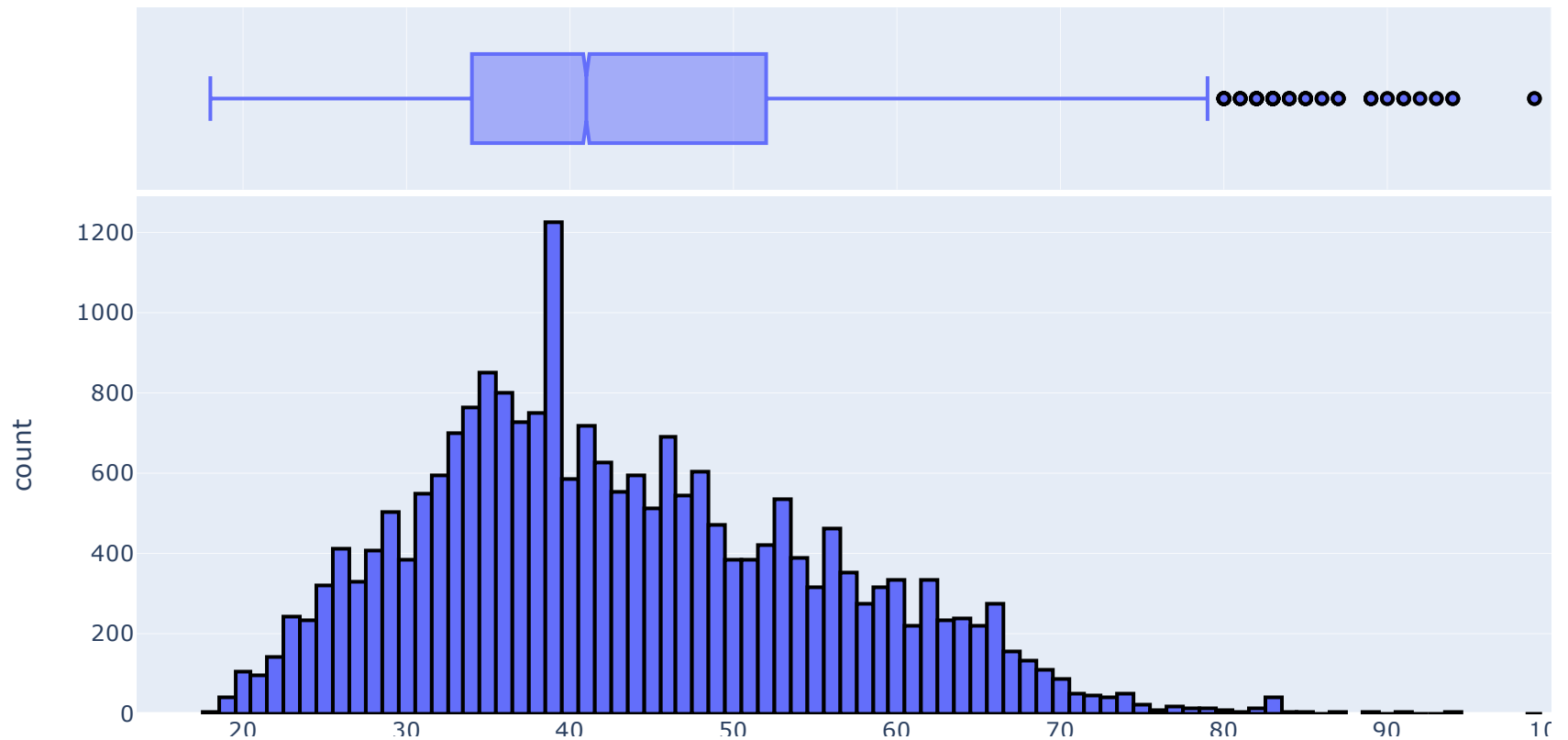


According to the graph above, frequency of the Rating 5 is pretty high compared to the others.

Distribution of the Age of the Customers

```
In [7]: 1 fig = px.histogram(data['Age'], marginal='box',
2               labels={'value': 'Age'})
3
4 fig.update_traces(marker=dict(line=dict(color='#000000', width=2)))
5 fig.update_layout(title_text='Distribution of the Age of the Customers',
6               title_x=0.5, title_font=dict(size=20))
7 fig.show()
```

Distribution of the Age of the Customers

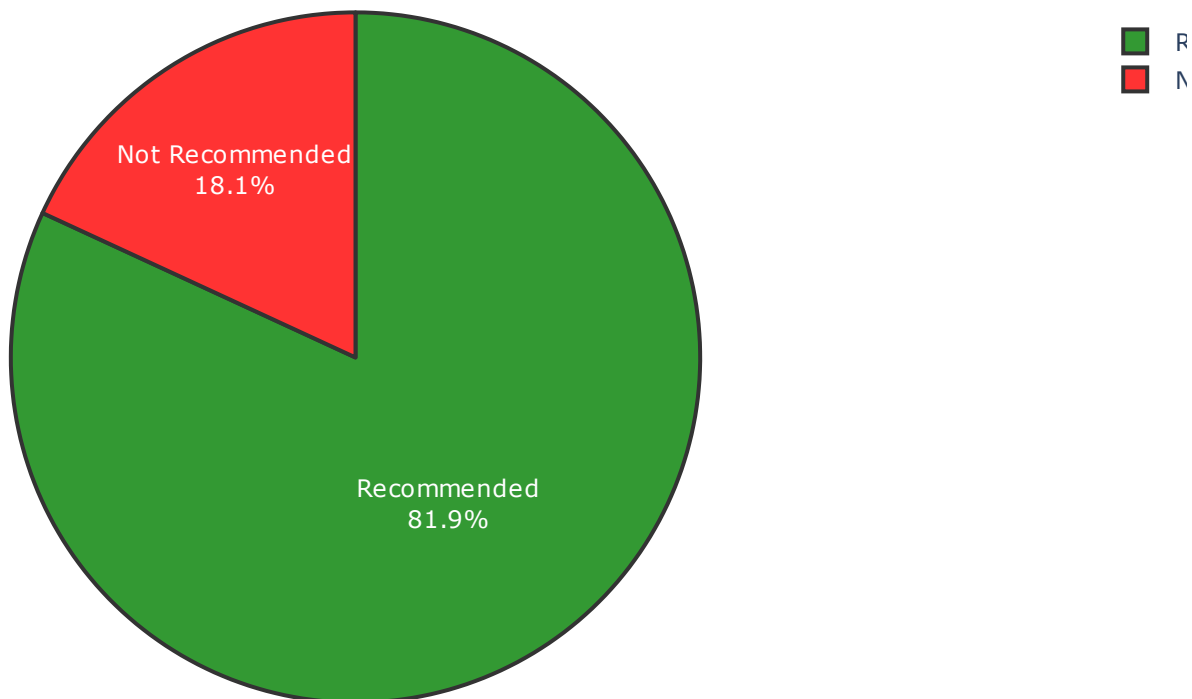


As you can see from the 'Distribution of the Age of the Customers' graph, the age of the customers is usually distributed between 34 and 52. We have outliers that customers older than 80.

Distribution of the Recommendations

```
In [8]: 1 labels = ['Recommended', 'Not Recommended']
2 values = [data[data['Recommended IND'] == 1]['Recommended IND'].value_counts()[1],
3           data[data['Recommended IND'] == 0]['Recommended IND'].value_counts()[0]]
4 colors = ['green', 'red']
5
6 fig = go.Figure(data=[go.Pie(labels=labels, values=values, opacity=0.8)])
7 fig.update_traces(textinfo='percent+label', marker=dict(line=dict(color='#000000', width=2), colors=col
8 fig.update_layout(title_text='Distribution of the Recommendations', title_x=0.5, title_font=dict(size=2
9 fig.show()
```

Distribution of the Recommendations

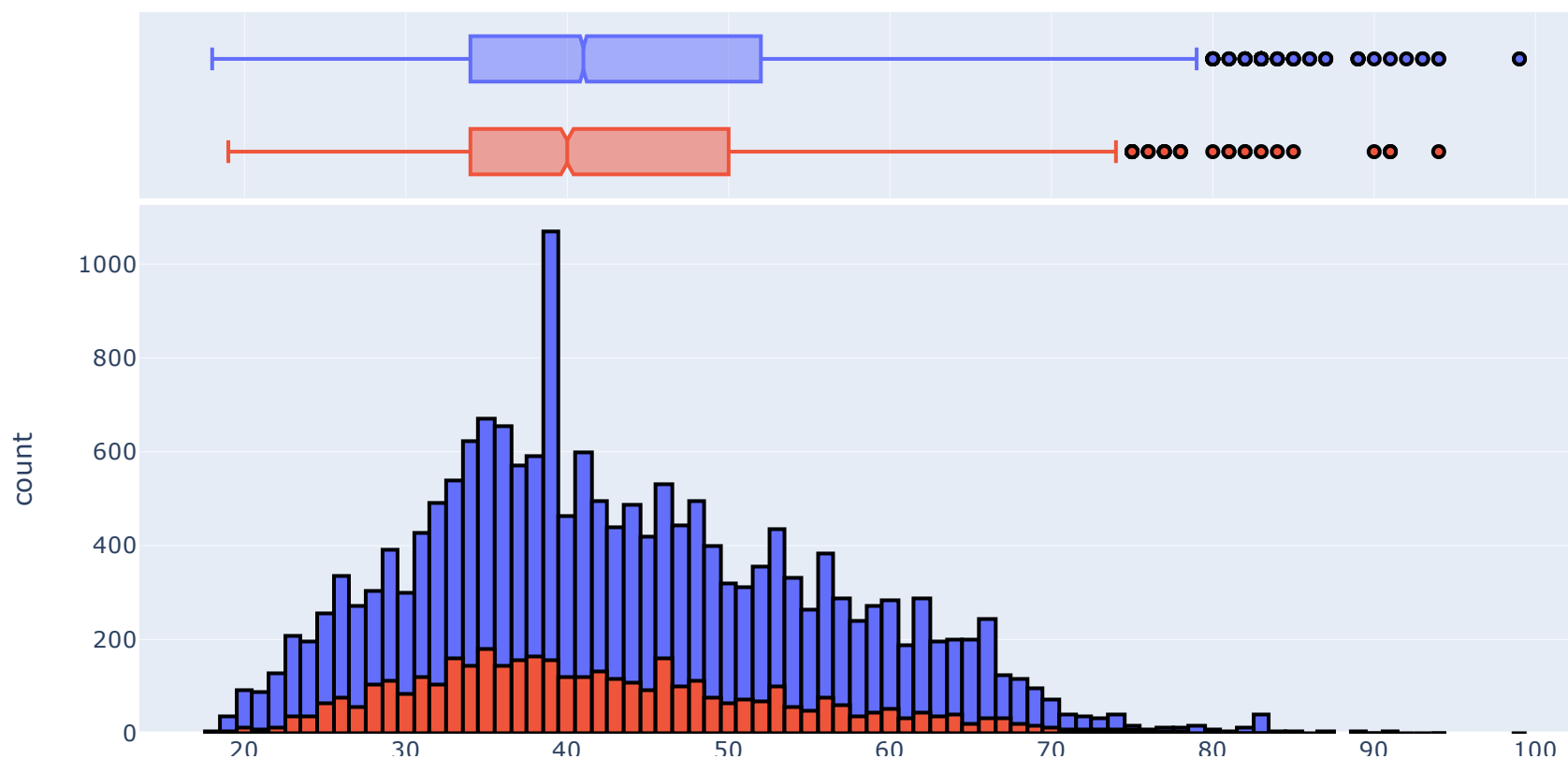


According to this pie chart, the most of the sales are Recommended.

Distribution of the Age and Recommendation

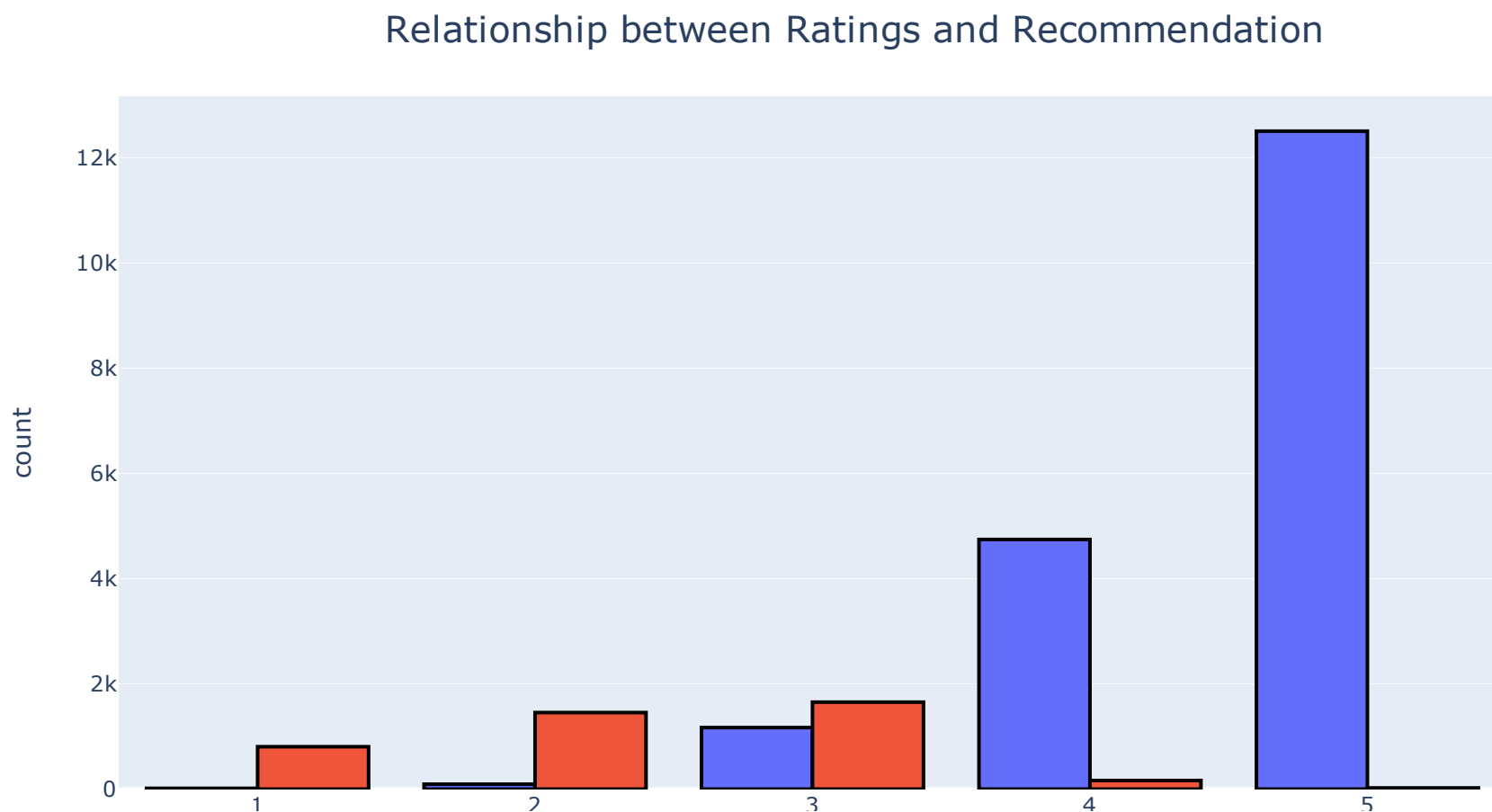
```
In [9]: 1 fig = px.histogram(data['Age'], color=data['Recommended IND'],
2                  labels={'value': 'Age',
3                        'color': 'Recommended'}, marginal='box')
4 fig.update_traces(marker=dict(line=dict(color='#000000', width=2)))
5 fig.update_layout(title_text='Distribution of the Age and Recommendation',
6                  title_x=0.5, title_font=dict(size=20))
7 fig.update_layout(barmode='overlay')
8 fig.show()
```

Distribution of the Age and Recommendation



Relationship between Ratings and Recommendation

```
In [10]: 1 fig = px.histogram(data['Rating'], color=data['Recommended IND'],
2                  labels={'value': 'Rating',
3                        'color': 'Recommended?'})
4 fig.update_layout(bargap=0.2)
5 fig.update_traces(marker=dict(line=dict(color='#000000', width=2)))
6 fig.update_layout(title_text='Relationship between Ratings and Recommendation',
7                  title_x=0.5, title_font=dict(size=20))
8 fig.update_layout(barmode='group')
9 fig.show()
```



According to this graph above, almost all the Rating 5 and Rating 4 data points are recommended.

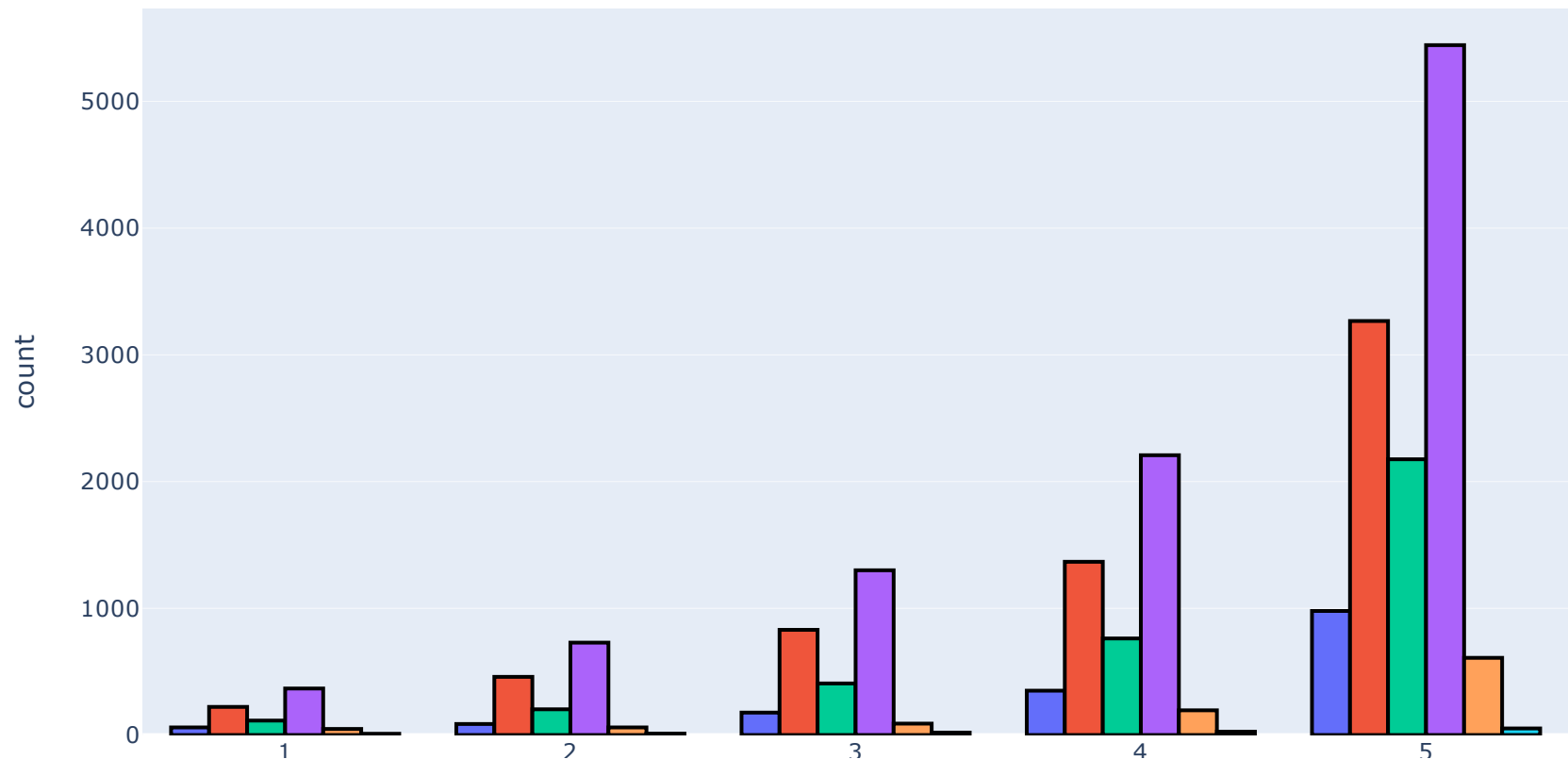
In addition, Rating 1 and Rating 2 data points have almost no recommendations.

For the further steps, I would create a common rating point with the Rating 4 and Rating 5 as well as Rating 1 and Rating 2. In this way, I would shrink the labels therefore, the model would perform better.

Relationship between Ratings and Departments

```
In [11]: 1 fig = px.histogram(data['Rating'], color=data['Department Name'],  
2                     labels={'value': 'Rating',  
3                           'color': 'Department Name'})  
4 fig.update_layout(bargap=0.2)  
5 fig.update_traces(marker=dict(line=dict(color='#000000', width=2)))  
6 fig.update_layout(title_text='Relationship between Ratings and Departments',  
7                   title_x=0.5, title_font=dict(size=20))  
8 fig.update_layout(barmode='group')  
9 fig.show()
```

Relationship between Ratings and Departments

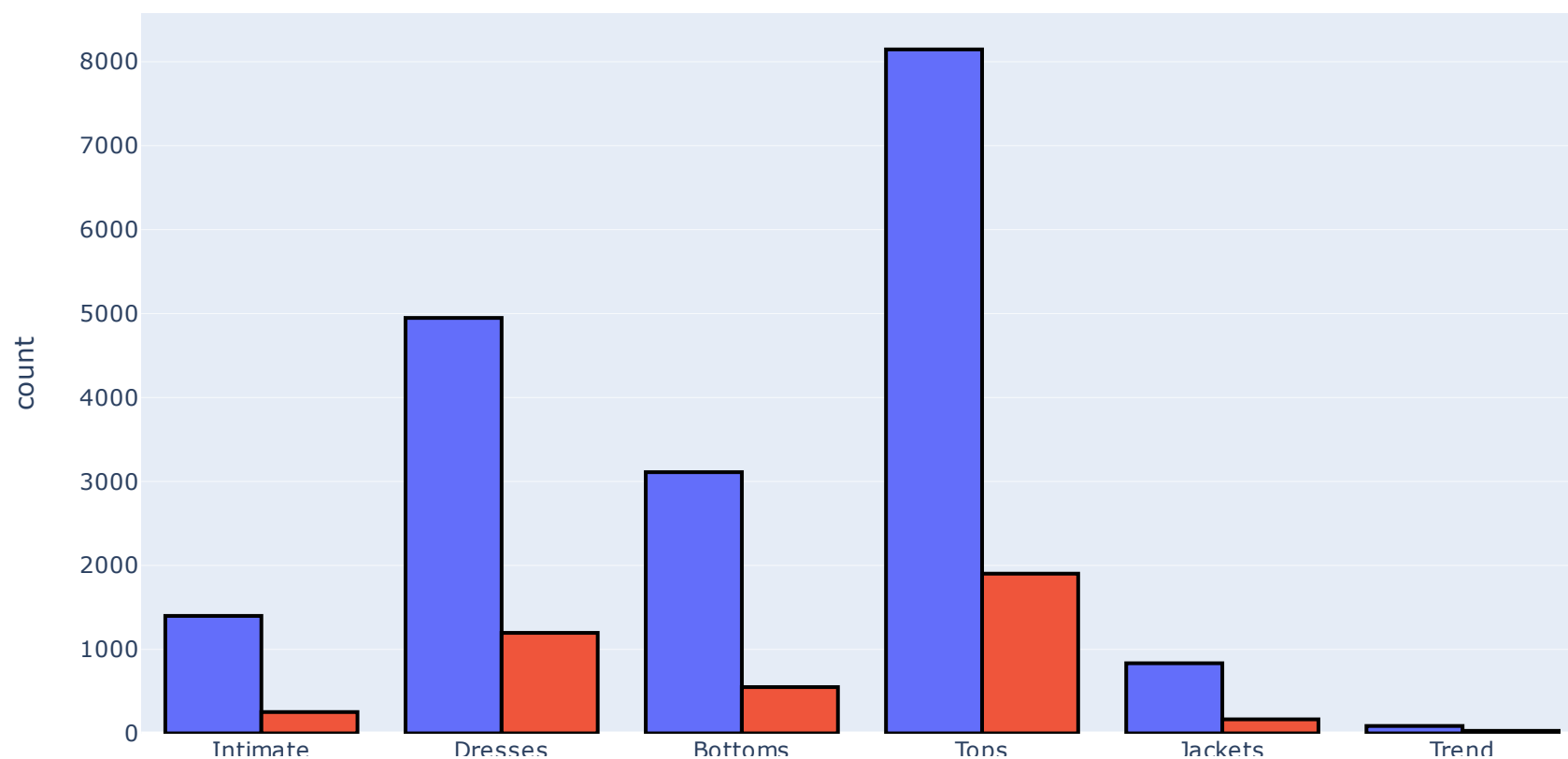


According to the graph above, Tops and Dresses have the most of the rating points. Trend and Jackets have the least.

Department and Recommendation Distribution

```
In [12]: 1 fig = px.histogram(data['Department Name'], color=data['Recommended IND'],  
2                    labels={'value': 'Department Name',  
3                    'color': 'Recommended?'})  
4 fig.update_layout(bargap=0.2)  
5 fig.update_traces(marker=dict(line=dict(color='#000000', width=2)))  
6 fig.update_layout(title_text='Department Name and Recommendation Distribution',  
7                    title_x=0.5, title_font=dict(size=20))  
8 fig.update_layout(barmode='group')  
9 fig.show()
```

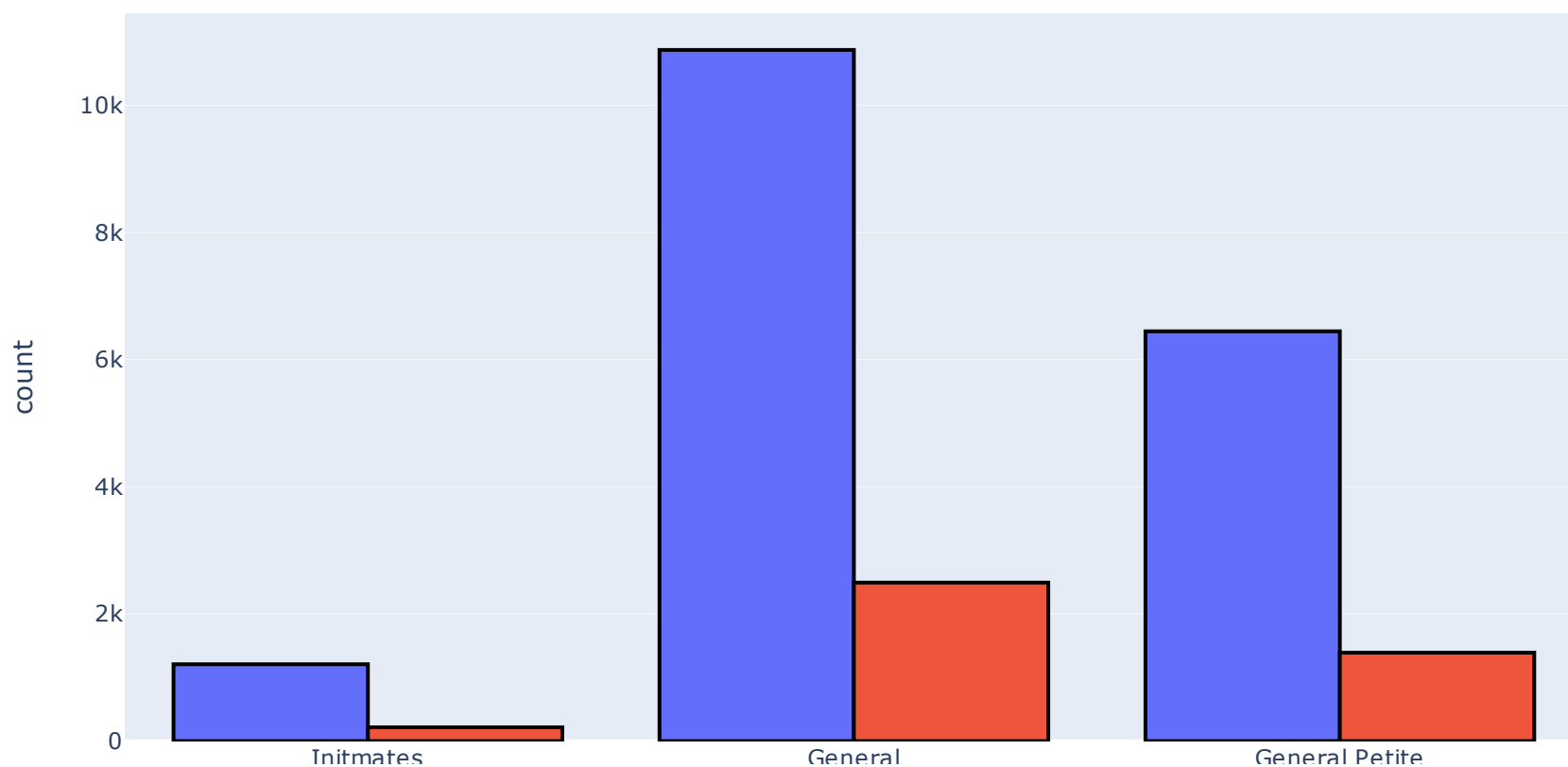
Department Name and Recommendation Distribution



Division and Recommendation Distribution

```
In [13]: 1 fig = px.histogram(data['Division Name'], color=data['Recommended IND'],
2                  labels={'value': 'Division Name',
3                        'color': 'Recommended?'})
4 fig.update_layout(bargap=0.2)
5 fig.update_traces(marker=dict(line=dict(color='#000000', width=2)))
6 fig.update_layout(title_text='Division Name and Recommendation Distribution',
7                  title_x=0.5, title_font=dict(size=20))
8 fig.update_layout(barmode='group')
9 fig.show()
```

Division Name and Recommendation Distribution



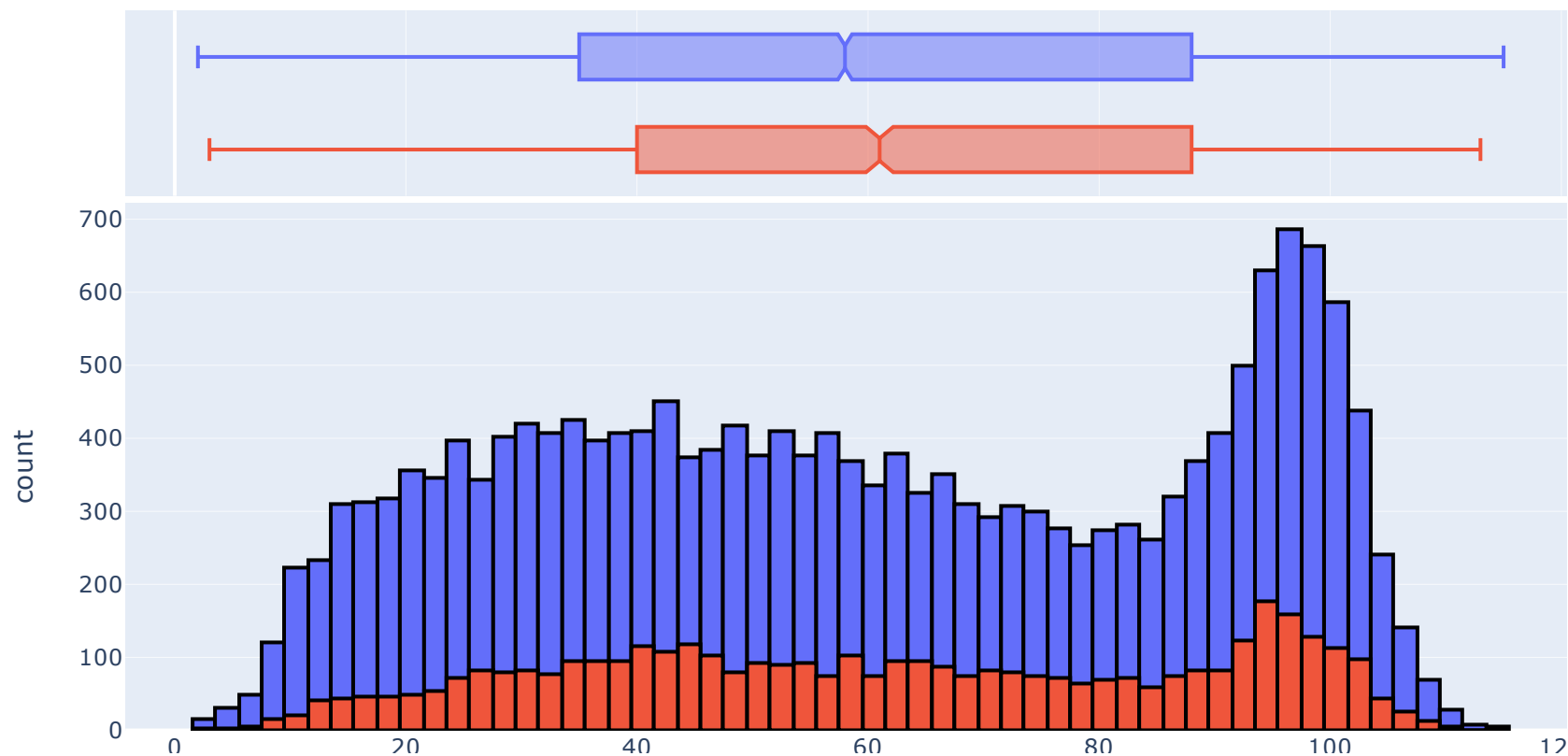
Distribution of the Length of the Texts

```

In [14]: 1 data['length_of_text'] = [len(i.split(' ')) for i in data['Review Text']]
2 fig = px.histogram(data['length_of_text'], marginal='box',
3                   labels={"value": "Length of the Text",
4                           "color": 'Recommended'},
5                   color=data['Recommended IND'])
6
7 fig.update_traces(marker=dict(line=dict(color='#000000', width=2)))
8 fig.update_layout(title_text='Distribution of the Length of the Texts',
9                   title_x=0.5, title_font=dict(size=20))
10 fig.update_layout(barmode='overlay')
11 fig.show()

```

Distribution of the Length of the Texts

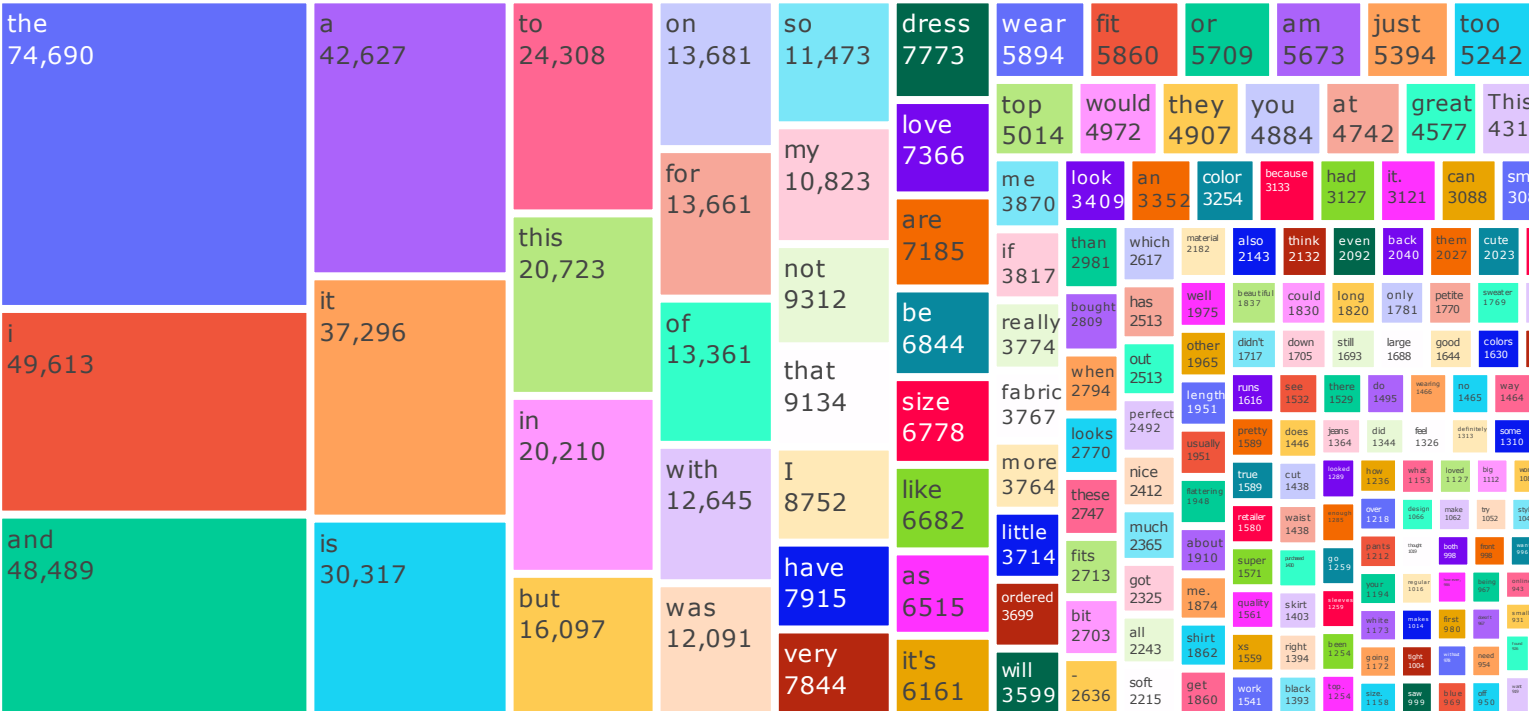


As you can see from the figure above, Recommended and not Recommended products almost have the same distribution length of text.

Top Frequent 200 Words in the Dataset (Before Cleaning)

```
1 FreqOfWords = data['Review Text'].str.split(expand=True).stack().value_counts()
2 FreqOfWords_top200 = FreqOfWords[:200]
3
4 fig = px.treemap(FreqOfWords_top200, path=[FreqOfWords_top200.index], values=0)
5 fig.update_layout(title_text='Top Frequent 200 Words in the Dataset (Before Cleaning)',
6                   title_x=0.5, title_font=dict(size=20)
7                   )
8 fig.update_traces(textinfo="label+value")
9 fig.show()
```

Top Frequent 200 Words in the Dataset (Before Cleaning)



According to this Treemap above, the top frequent 200 words usually include a stopword. For the further step of this notebook, I will remove them from the text.

Data Preprocessing

```
In [16]: 1 # Lower Character all the Texts
          2 data['Review Text'] = data['Review Text'].str.lower()
          3 data['Review Text'].head()
```

```
Out[16]: 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...
          Name: Review Text, dtype: object
```

```
In [17]: 1 # Removing Punctuations and Numbers from the Text
          2 def remove_punctuations_numbers(inputs):
          3     return re.sub(r'^a-zA-Z', ' ', inputs)
          4
          5
          6 data['Review Text'] = data['Review Text'].apply(remove_punctuations_numbers)
```

Now we will remove all punctuations and numbers from the all dataframe. They will be not usefull for my model training.

Tokenizing with NLTK

```
In [18]: 1 def tokenization(inputs): # Ref.1
          2     return word_tokenize(inputs)
          3
          4
          5 data['text_tokenized'] = data['Review Text'].apply(tokenization)
          6 data['text_tokenized'].head()
```

```
Out[18]: 0 [absolutely, wonderful, silky, and, sexy, and,...
          1 [love, this, dress, it, s, sooo, pretty, i, ha...
          2 [i, had, such, high, hopes, for, this, dress, ...
          3 [i, love, love, love, this, jumpsuit, it, s, f...
          4 [this, shirt, is, very, flattering, to, all, d...
          Name: text_tokenized, dtype: object
```

Stopwords Removal

```
In [20]: 1 import nltk
          2 nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\NISHANTH\AppData\Roaming\nltk_data...
[nltk_data] Unzipping corpora\stopwords.zip.
```

```
Out[20]: True
```

```
In [21]: 1 stop_words = set(stopwords.words('english'))
2 stop_words.remove('not')
3
4
5 def stopwords_remove(inputs): # Ref.2
6     return [k for k in inputs if k not in stop_words]
7
8
9 data['text_stop'] = data['text_tokenized'].apply(stopwords_remove)
10 data['text_stop'].head()
```

```
Out[21]: 0 [absolutely, wonderful, silky, sexy, comfortable]
1 [love, dress, sooo, pretty, happened, find, st...
2 [high, hopes, dress, really, wanted, work, ini...
3 [love, love, love, jumpsuit, fun, flirty, fabu...
4 [shirt, flattering, due, adjustable, front, ti...
Name: text_stop, dtype: object
```

Lemmatization

```
In [22]: 1 lemmatizer = WordNetLemmatizer()
2
3
4 def lemmatization(inputs): # Ref.1
5     return [lemmatizer.lemmatize(word=kk, pos='v') for kk in inputs]
6
7
8 data['text_lemmatized'] = data['text_stop'].apply(lemmatization)
9 data['text_lemmatized'].head()
```

```
Out[22]: 0 [absolutely, wonderful, silky, sexy, comfortable]
1 [love, dress, sooo, pretty, happen, find, stor...
2 [high, hop, dress, really, want, work, initial...
3 [love, love, love, jumpsuit, fun, flirty, fabu...
4 [shirt, flatter, due, adjustable, front, tie, ...
Name: text_lemmatized, dtype: object
```

```
In [23]: 1 # Removing Words less than length 2
2 def remove_less_than_2(inputs): # Ref.1
3     return [j for j in inputs if len(j) > 2]
4
5
6 data['final'] = data['text_lemmatized'].apply(remove_less_than_2)
```

```
In [24]: 1 # Joining Tokens into Sentences
2 data['final'] = data['final'].str.join(' ')
3 data['final'].head()
```

```
Out[24]: 0      absolutely wonderful silky sexy comfortable
1      love dress sooo pretty happen find store glad ...
2      high hop dress really want work initially orde...
3      love love love jumpsuit fun flirty fabulous ev...
4      shirt flatter due adjustable front tie perfect...
Name: final, dtype: object
```

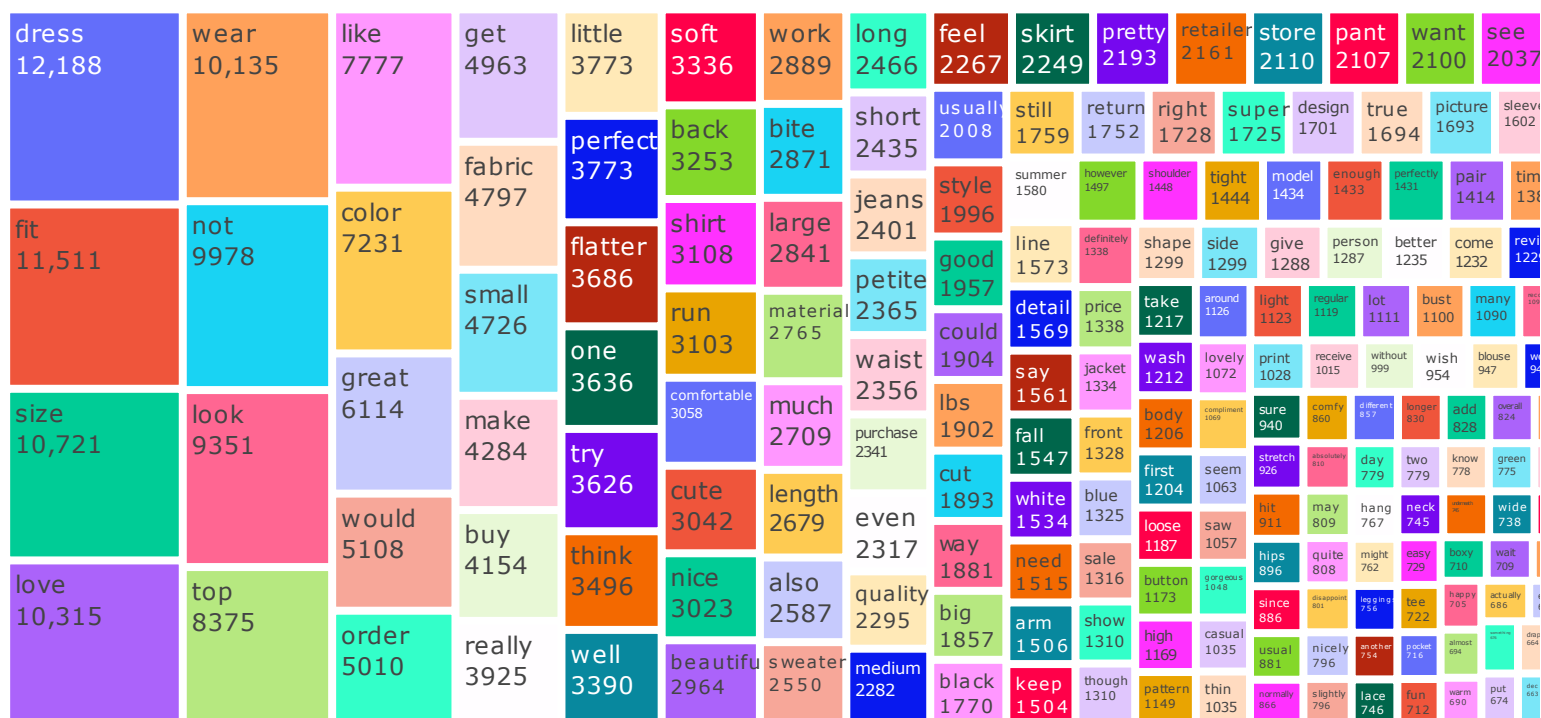
Top Frequent 200 Words in the Dataset (After Cleaning)

```

In [25]: 1 FreqOfWords = data['final'].str.split(expand=True).stack().value_counts()
2 FreqOfWords_top200 = FreqOfWords[:200]
3
4 fig = px.treemap(FreqOfWords_top200, path=[FreqOfWords_top200.index], values=0)
5 fig.update_layout(title_text='Top Frequent 200 Words in the Dataset (After Cleaning)',
6                   title_x=0.5, title_font=dict(size=20)
7                   )
8 fig.update_traces(textinfo="label+value")
9 fig.show()

```

Top Frequent 200 Words in the Dataset (After Cleaning)



As you can see from the Treemap above, all of the words are unique words and there are no stopwords in this set. Most words are 'dress', 'fit' and 'size'. Due to we are dealing with the clothing review dataset, this is pretty reasonable.

WordCloud of the Recommended Reviews

```
In [26]: 1 data_recommended = data[data['Recommended IND'] == 1] # Dataframe that only includes recommended reviews
2 data_not_recommended = data[data['Recommended IND'] == 0] # Dataframe that only includes not recommended reviews
3
4 WordCloud_recommended = WordCloud(max_words=500,
5                                   random_state=30,
6                                   collocations=True).generate(str((data_recommended['final'])))
7
8 plt.figure(figsize=(15, 8))
9 plt.imshow(WordCloud_recommended, interpolation='bilinear')
10 plt.title('WordCloud of the Recommended Reviews', fontsize=20)
11 plt.axis("off")
12 plt.show()
```

```
In [27]: 1 FreqOfWords = data_recommended['final'].str.split(expand=True).stack().value_counts()
2 FreqOfWords_top200 = FreqOfWords[:200]
3
4 fig = px.treemap(FreqOfWords_top200, path=[FreqOfWords_top200.index], values=0)
5 fig.update_layout(title_text='Top Frequent 200 Words in the Recommended Reviews',
6                   title_x=0.5, title_font=dict(size=20)
7                   )
8 fig.update_traces(textinfo="label+value")
9 fig.show()
```

Top Frequent 200 Words in the Recommended Reviews



WordCloud of the Not Recommended Reviews

```
In [28]: 1 WordCloud_not_recommended = WordCloud(max_words=500,  
2                                           random_state=30,  
3                                           collocations=True).generate(str((data_not_recommended['final'])))  
4  
5 plt.figure(figsize=(15, 8))  
6 plt.imshow(WordCloud_not_recommended, interpolation='bilinear')  
7 plt.title('WordCloud of the Not Recommended Reviews', fontsize=20)  
8 plt.axis("off")  
9 plt.show()
```

WordCloud of the Not Recommended Reviews



```
In [29]: 1 FreqOfWords = data_not_recommended['final'].str.split(expand=True).stack().value_counts()
2 FreqOfWords_top200 = FreqOfWords[:200]
3
4 fig = px.treemap(FreqOfWords_top200, path=[FreqOfWords_top200.index], values=0)
5 fig.update_layout(title_text='Top Frequent 200 Words in the Not Recommended Reviews',
6                   title_x=0.5, title_font=dict(size=20)
7                   )
8 fig.update_traces(textinfo="label+value")
9 fig.show()
```

Top Frequent 200 Words in the Not Recommended Reviews



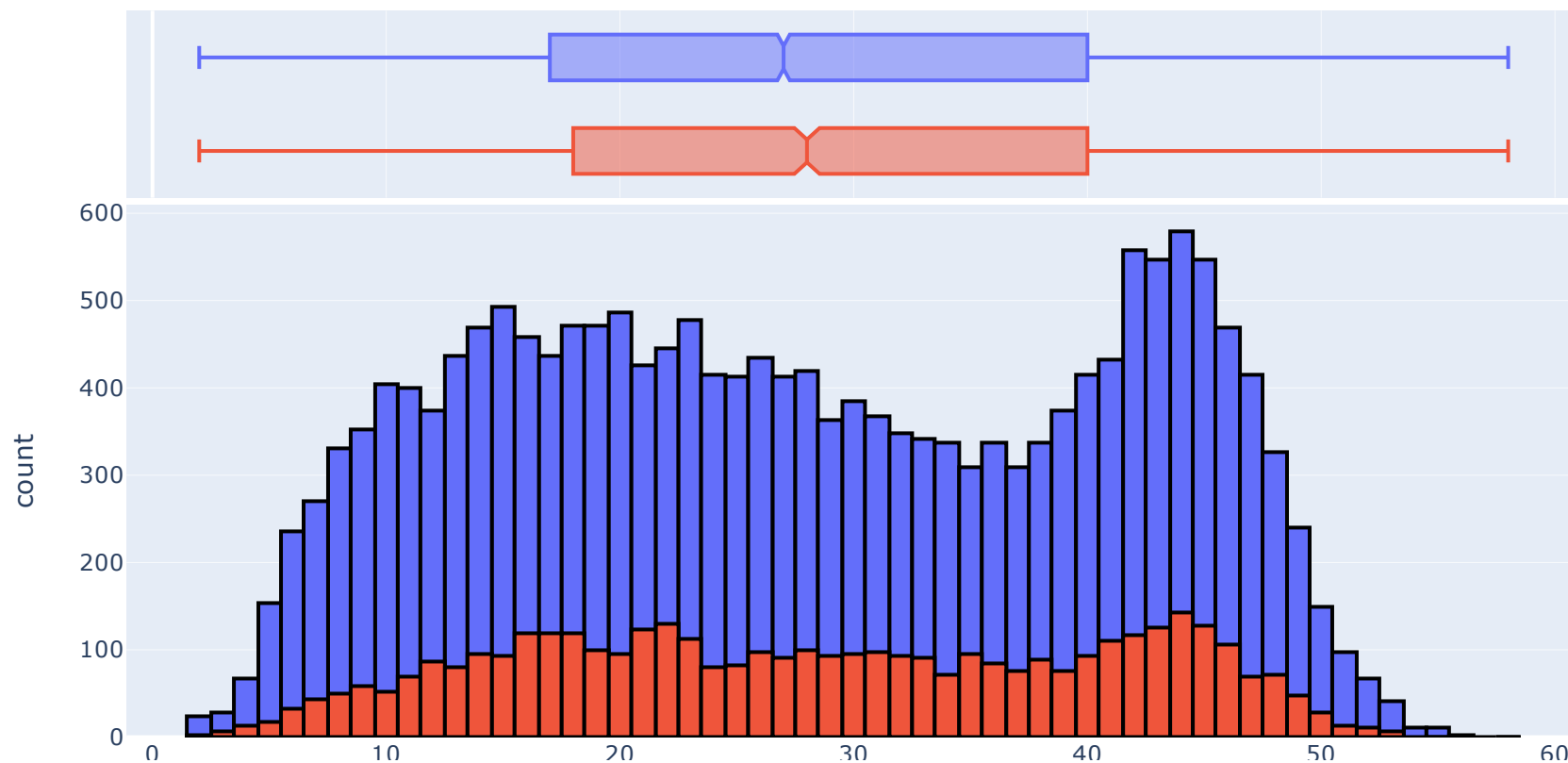
Distribution of the Length of the Texts after Cleaning

```

In [30]: 1 data['length_of_text'] = [len(i.split(' ')) for i in data['final']]
2 fig = px.histogram(data['length_of_text'], marginal='box',
3                   labels={"value": "Length of the Text",
4                           "color": 'Recommended?'},
5                   color=data['Recommended IND'])
6
7 fig.update_traces(marker=dict(line=dict(color='#000000', width=2)))
8 fig.update_layout(title_text='Distribution of the Length of the Texts after Cleaning',
9                   title_x=0.5, title_font=dict(size=20))
10 fig.update_layout(barmode='overlay')
11 fig.show()

```

Distribution of the Length of the Texts after Cleaning



Train-Test-Validation Split

```
In [31]: 1 # I will only use Text data to predict Recommendation
        2 y = data['Recommended IND']
        3 X = data['final']
        4
        5 X.head()
```

```
Out[31]: 0      absolutely wonderful silky sexy comfortable
        1      love dress sooo pretty happen find store glad ...
        2      high hop dress really want work initially orde...
        3      love love love jumpsuit fun flirty fabulous ev...
        4      shirt flatter due adjustable front tie perfect...
        Name: final, dtype: object
```

```
In [32]: 1 # Train-Test-Validation Split
        2 x, X_test, y, y_test = train_test_split(X, y, test_size=0.2, random_state=13) # Test: %20
        3
        4 X_train, X_val, y_train, y_val = train_test_split(x, y, test_size=0.25, random_state=13) # Val: %20
        5
        6 print('Shape of the X_train:', X_train.shape)
        7 print('Shape of the X_test:', X_test.shape)
        8 print('Shape of the X_val:', X_val.shape)
        9 print('--'*20)
       10 print('Shape of the y_train:', y_train.shape)
       11 print('Shape of the y_test:', y_test.shape)
       12 print('Shape of the y_val:', y_val.shape)
```

```
Shape of the X_train: (13576,)
Shape of the X_test: (4526,)
Shape of the X_val: (4526,)
```

```
-----
```

```
Shape of the y_train: (13576,)
Shape of the y_test: (4526,)
Shape of the y_val: (4526,)
```

Tokenizing with Tensorflow

```
In [33]: 1 num_words = 10000
2 tokenizer = Tokenizer(num_words=num_words, oov_token='<OOV>')
3 tokenizer.fit_on_texts(X_train)
4
5 Tokenized_train = tokenizer.texts_to_sequences(X_train)
6 Tokenized_val = tokenizer.texts_to_sequences(X_val)
7
8 print('Non-tokenized Version: ', X_train[0])
9 print('Tokenized Version: ', tokenizer.texts_to_sequences([X_train[0]]))
10 print('---'*20)
11 print('Non-tokenized Version: ', X_train[80])
12 print('Tokenized Version: ', tokenizer.texts_to_sequences([X_train[80]]))
```

Non-tokenized Version: absolutely wonderful silky sexy comfortable

Tokenized Version: [[161, 366, 748, 445, 33]]

Non-tokenized Version: usually petite since dress not come petites try fit lbs dress hit knee hem bite no
t overwhelm dress look stun great vibrant color dark hair make classic elegant dress look contemporary sty
lish try store salesperson others happen see rave tell grab glad plan wear spring daughte

Tokenized Version: [[61, 47, 150, 2, 7, 109, 769, 23, 3, 68, 2, 146, 269, 223, 38, 7, 746, 2, 8, 397, 12,
356, 11, 278, 1025, 18, 342, 459, 2, 8, 2151, 344, 23, 57, 1885, 350, 601, 62, 1723, 432, 631, 218, 313,
6, 212, 6315]]

Padding the Datasets

```
In [34]: 1 maxlen = 50
2 Padded_train = pad_sequences(Tokenized_train, maxlen=maxlen, padding='pre')
3 Padded_val = pad_sequences(Tokenized_val, maxlen=maxlen, padding='pre')
```

ANN Model Creation

```
In [35]: 1 # Creating the Model
2 model = Sequential()
3
4 model.add(Embedding(num_words, 16, input_length=maxlen))
5 model.add(Dropout(0.2))
6
7 model.add(GlobalAvgPool1D())
8 model.add(Dropout(0.5))
9
10 model.add(Dense(1, activation='sigmoid'))
11
12 opt = tf.optimizers.Adam(lr=0.55e-3) # Learning Rate
13
14 model.compile(loss='binary_crossentropy', optimizer=opt, metrics=['accuracy'])
15
16 model.summary()
```

Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
embedding (Embedding)	(None, 50, 16)	160000
dropout (Dropout)	(None, 50, 16)	0
global_average_pooling1d (Gl	(None, 16)	0
dropout_1 (Dropout)	(None, 16)	0
dense (Dense)	(None, 1)	17
=====		
Total params: 160,017		
Trainable params: 160,017		
Non-trainable params: 0		

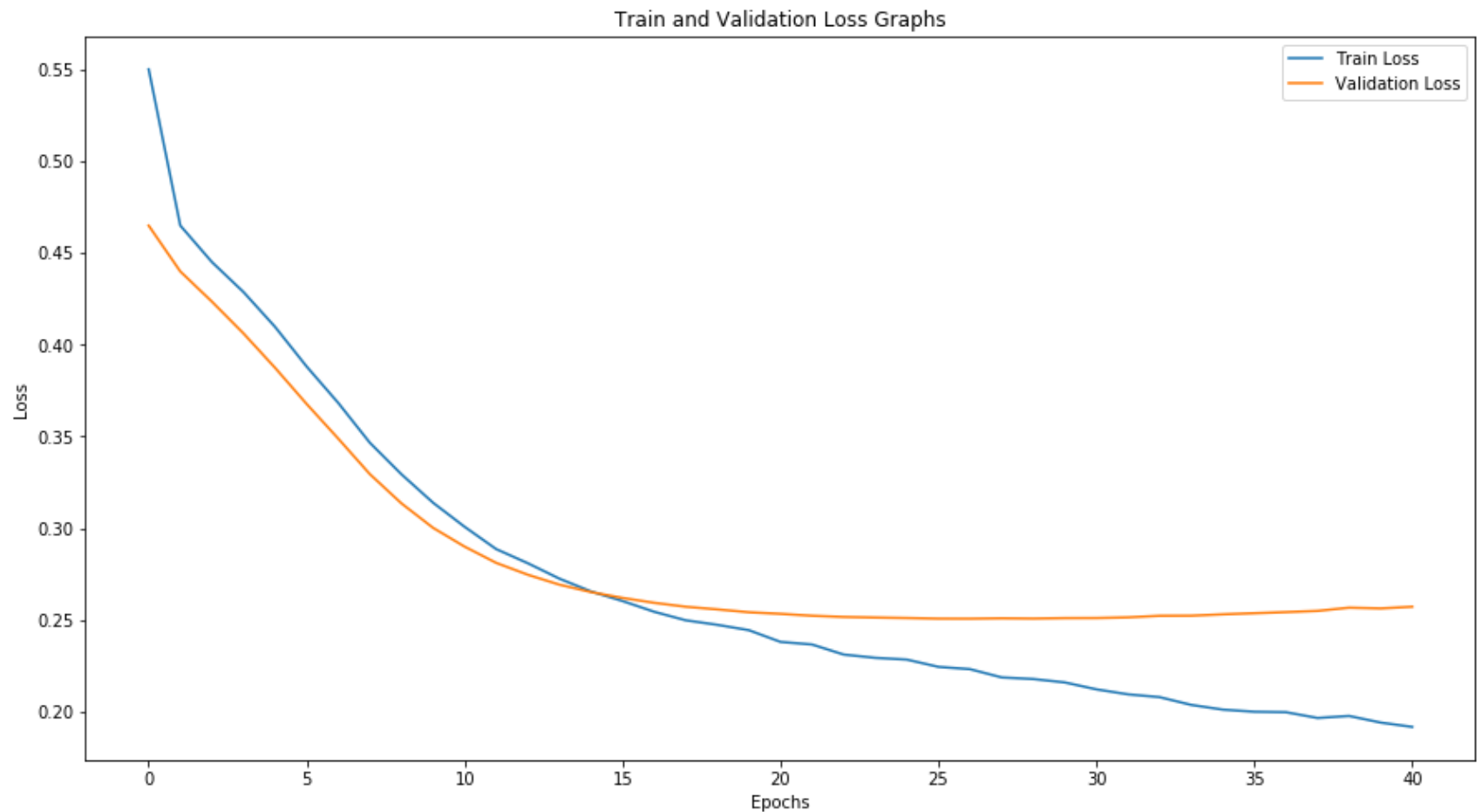
```
In [36]: 1 # Training the Model
2 early_stopping = tf.keras.callbacks.EarlyStopping(monitor='val_accuracy', mode='auto', patience=5,
3                                                    restore_best_weights=True)
4
5 epochs = 100
6 hist = model.fit(Padded_train, y_train, epochs=epochs,
7                 validation_data=(Padded_val, y_val),
8                 callbacks=[early_stopping], batch_size=32)
```

```
Epoch 1/100
425/425 [=====] - 6s 7ms/step - loss: 0.6074 - accuracy: 0.8130 - val_loss: 0.
4649 - val_accuracy: 0.8239
Epoch 2/100
425/425 [=====] - 2s 5ms/step - loss: 0.4773 - accuracy: 0.8132 - val_loss: 0.
4400 - val_accuracy: 0.8239
Epoch 3/100
425/425 [=====] - 2s 5ms/step - loss: 0.4556 - accuracy: 0.8144 - val_loss: 0.
4235 - val_accuracy: 0.8239
Epoch 4/100
425/425 [=====] - 2s 5ms/step - loss: 0.4340 - accuracy: 0.8194 - val_loss: 0.
4062 - val_accuracy: 0.8239
Epoch 5/100
425/425 [=====] - 2s 5ms/step - loss: 0.4161 - accuracy: 0.8183 - val_loss: 0.
3875 - val_accuracy: 0.8243
Epoch 6/100
425/425 [=====] - 2s 5ms/step - loss: 0.3985 - accuracy: 0.8161 - val_loss: 0.
3677 - val_accuracy: 0.8261
Epoch 7/100
425/425 [=====] - 2s 5ms/step - loss: 0.3678 - accuracy: 0.8274 - val_loss: 0.
```

Train and Validation Loss Graphs


```
In [38]: 1 plt.figure(figsize=(15, 8))
2         plt.plot(hist.history['loss'], label='Train Loss')
3         plt.plot(hist.history['val_loss'], label='Validation Loss')
4         plt.title('Train and Validation Loss Graphs')
5         plt.xlabel('Epochs')
6         plt.ylabel('Loss')
7         plt.legend()
```

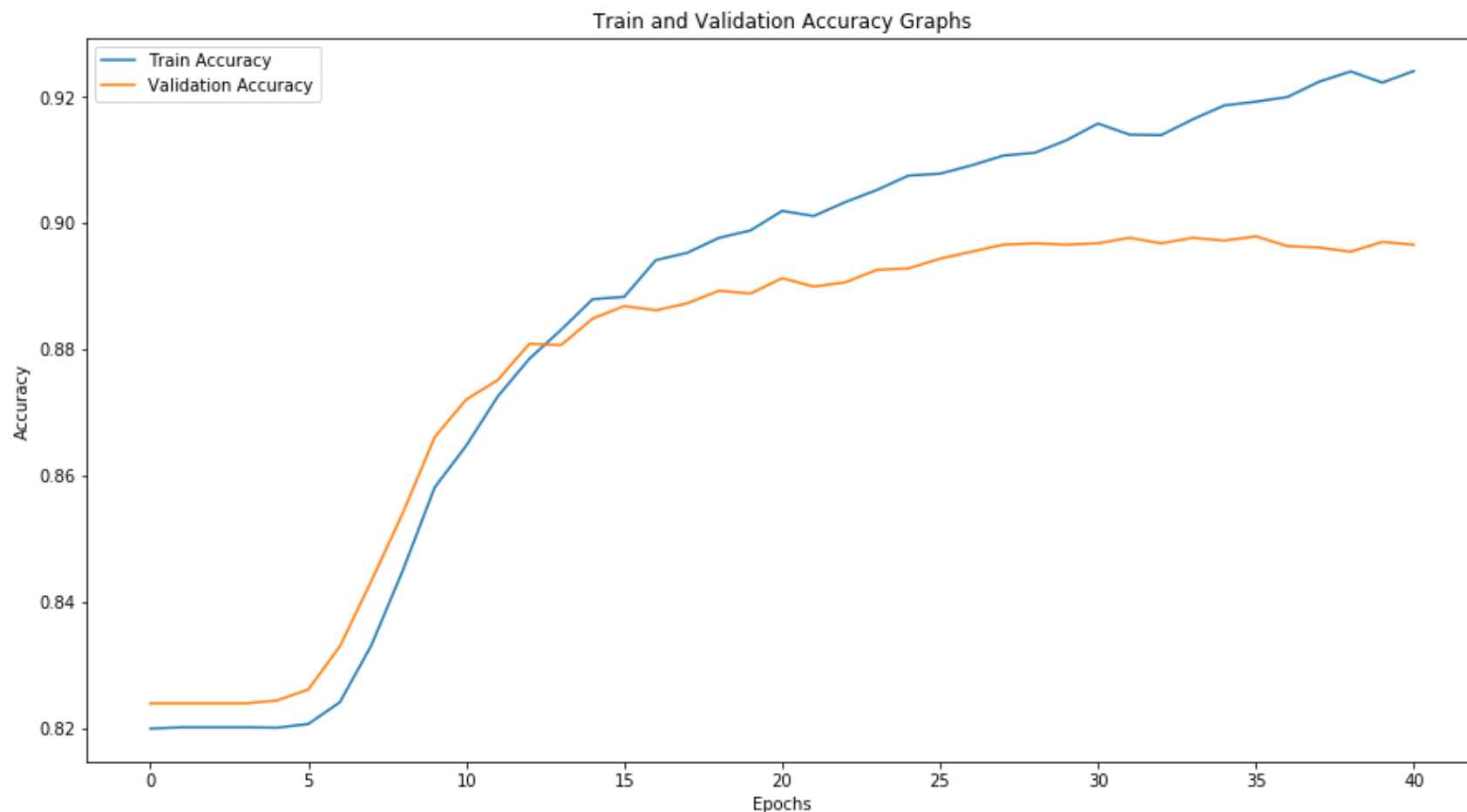
Out[38]: <matplotlib.legend.Legend at 0x15f24dae448>



Train and Validation Accuracy Graphs

```
In [39]: 1 plt.figure(figsize=(15, 8))
2 plt.plot(hist.history['accuracy'], label='Train Accuracy')
3 plt.plot(hist.history['val_accuracy'], label='Validation Accuracy')
4 plt.title('Train and Validation Accuracy Graphs')
5 plt.xlabel('Epochs')
6 plt.ylabel('Accuracy')
7 plt.legend()
```

Out[39]: <matplotlib.legend.Legend at 0x15f254c6ec8>



Preparing the Test Data

```
In [40]: 1 X_test = X_test.apply(tokenization)
2 X_test = X_test.apply(stopwords_remove)
3 X_test = X_test.apply(lemmatization)
4 X_test = X_test.str.join(' ')
5
6 X_test.head()
```

```
Out[40]: 10818          low waisted weird liner
779      shirt not good look gal hips fit top tight ord...
10907     love dress long enough dramatic graze feet wit...
17442    understand pencil skirt gon body hug however r...
832      order shirt wear pair pant return one reason t...
Name: final, dtype: object
```

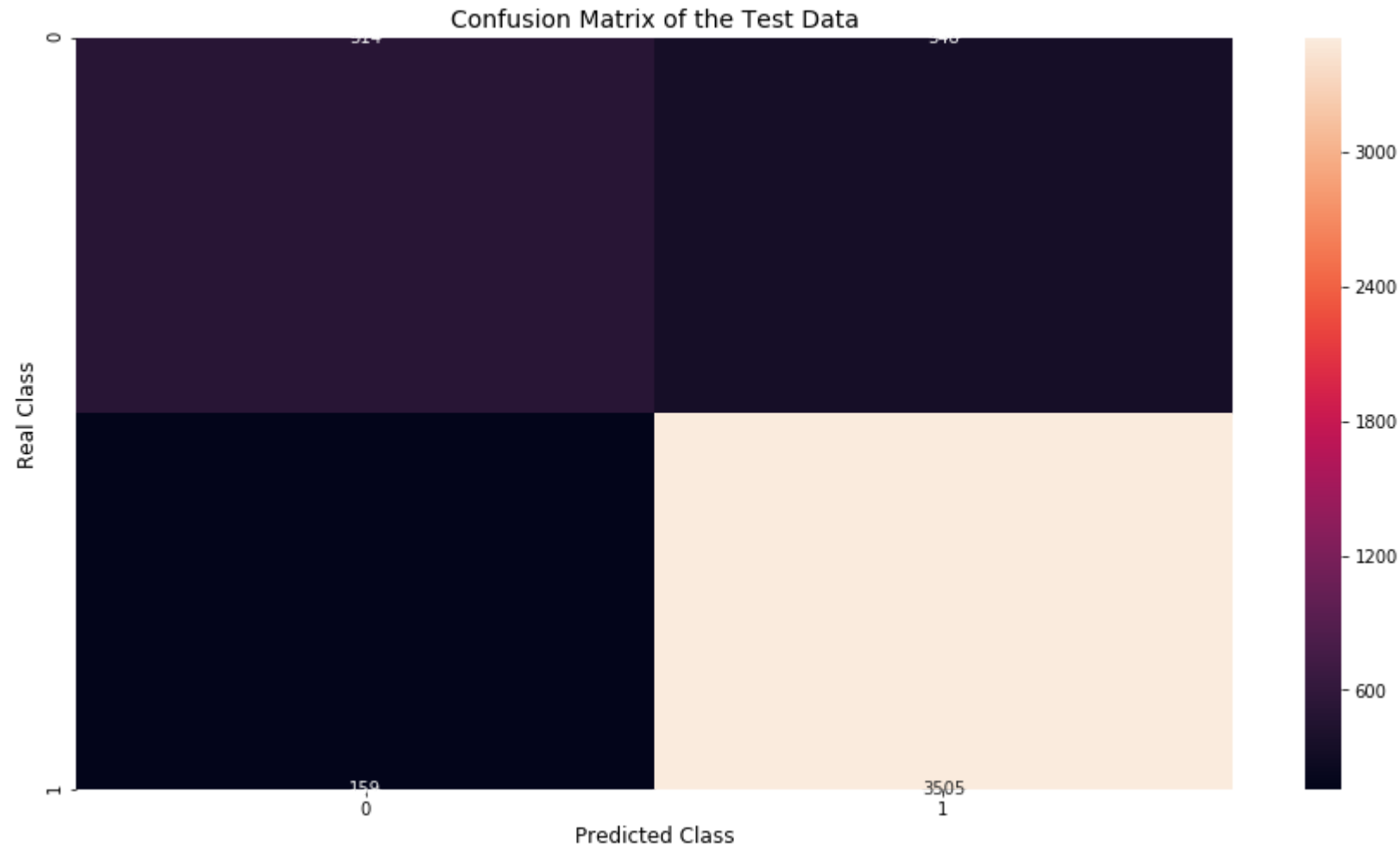
```
In [41]: 1 Tokenized_test = tokenizer.texts_to_sequences(X_test)
2 Padded_test = pad_sequences(Tokenized_test, maxlen=maxlen, padding='pre')
3
4 test_evaluate = model.evaluate(Padded_test, y_test)
```

```
142/142 [=====] - 0s 1ms/step - loss: 0.2700 - accuracy: 0.8889
```

Confusion Matrix of the Test Data

```
In [42]: 1 pred_train_lstm = model.predict(Padded_train)
2 pred_test_lstm = model.predict(Padded_test)
```

```
In [43]: 1 for i, x in enumerate(pred_test_lstm):
2         if 0 <= x < 0.49:
3             pred_test_lstm[i] = 0
4         else:
5             pred_test_lstm[i] = 1
6
7     for i, x in enumerate(pred_train_lstm):
8         if 0 <= x < 0.49:
9             pred_train_lstm[i] = 0
10        else:
11            pred_train_lstm[i] = 1
12
13    conf_mat = confusion_matrix(y_true=y_test, y_pred=pred_test_lstm)
14    plt.figure(figsize=(15, 8))
15    sns.heatmap(conf_mat, annot=True, fmt='g')
16    plt.title('Confusion Matrix of the Test Data', fontsize=14)
17    plt.ylabel('Real Class', fontsize=12)
18    plt.xlabel('Predicted Class', fontsize=12)
19    plt.show()
```



Evaluation Metrics of the LSTM Model

In [44]:

```
1 # Accuracy
2 train_acc_lstm = round(accuracy_score(y_train, pred_train_lstm) * 100, 2)
3 print('Train Accuracy of the LSTM: %', train_acc_lstm)
4 test_acc_lstm = round(accuracy_score(y_test, pred_test_lstm) * 100, 2)
5 print('Test Accuracy of the LSTM: %', test_acc_lstm)
6 print('--' * 20)
7
8 # Precision
9 train_precision_lstm = round(precision_score(y_train, pred_train_lstm) * 100, 2)
10 print('Train Precision of the LSTM: %', train_precision_lstm)
11 precision_lstm = round(precision_score(y_test, pred_test_lstm) * 100, 2)
12 print('Test Precision of the LSTM: %', precision_lstm)
13 print('--' * 20)
14
15 # Recall
16 train_recall_lstm = round(recall_score(y_train, pred_train_lstm) * 100, 2)
17 print('Train Recall of the LSTM: %', train_recall_lstm)
18 recall_lstm = round(recall_score(y_test, pred_test_lstm) * 100, 2)
19 print('Test Recall of the LSTM: %', recall_lstm)
```

Train Accuracy of the LSTM: % 92.35

Test Accuracy of the LSTM: % 88.8

Train Precision of the LSTM: % 93.94

Test Precision of the LSTM: % 90.97

Train Recall of the LSTM: % 96.93

Test Recall of the LSTM: % 95.66

Having Fun with the LSTM Model

```
In [45]: 1 def predict_recommendation(input_text): # The function for doing all the previous steps
2         input_text = input_text.lower()
3         input_text = re.sub(r'^a-zA-Z', ' ', input_text)
4         input_text = tokenization(input_text)
5         input_text = stopwords_remove(input_text)
6         input_text = lemmatization(input_text)
7         input_text = ' '.join(input_text)
8         input_text = tokenizer.texts_to_sequences([input_text])
9         input_text = pad_sequences(input_text, maxlen=maxlen, padding='pre')
10        input_text = model.predict(input_text)
11        if input_text >= 0.5:
12            input_text = f'Recommended with %{round(float(input_text*100), 2)}%'
13        else:
14            input_text = f'Not Recommended with %{round(float(input_text*100), 2)}%'
15
16        return print(input_text)
```

```
In [46]: 1 # This reviews above are taken from several websites for testing the model with real world data. You ca
2         predict_recommendation("The clothes are such poor quality and look nothing like they do on the website.
    <img alt="A horizontal scrollbar with a grey track and a white slider, indicating the text is scrollable." data-bbox="175 495 955 515"/>
Not Recommended with %16.57
```

```
In [47]: 1 predict_recommendation("Beautiful colour of lemon great fit and length here in three days all I need is
    <img alt="A horizontal scrollbar with a grey track and a white slider, indicating the text is scrollable." data-bbox="175 610 955 630"/>
Recommended with %94.93
```

```
In [48]: 1 predict_recommendation("As usual the clothes I ordered arrived quickly and were all a good fit, except
    <img alt="A horizontal scrollbar with a grey track and a white slider, indicating the text is scrollable." data-bbox="175 725 955 745"/>
Recommended with %83.63
```

```
In [49]: 1 predict_recommendation("I should've checked reviews before ordering... each item they sent was much wor
    <img alt="A horizontal scrollbar with a grey track and a white slider, indicating the text is scrollable." data-bbox="175 840 955 860"/>
Not Recommended with %2.59
```


In [50]: 1 predict_recommendation("cheap material that falls apart in seconds. Clothes look nothing like the picture")
Not Recommended with %19.39

In [51]: 1 predict_recommendation("Very fast dispatch and delivery. Clothes are always a consistent fit, good quality")
Recommended with %89.66

In [52]: 1 predict_recommendation("I have no complaints whatsoever, from ordering to getting my goods were excellent")
Recommended with %94.01

In [53]: 1 predict_recommendation("My dress had blue ink and biro stains on which was a real shame. I needed it for a while")
Not Recommended with %41.35

In [54]: 1 # Ref. 6 from now on
2 predict_recommendation("Sizes varied despite allegedly being the same size. Some of the quality was poor")
Not Recommended with %14.46

In [55]: 1 predict_recommendation("I do really like yours clothing, just find the sizing is slightly off, a normal fit")
Recommended with %89.15

In [56]: 1 predict_recommendation("I really love this dress. I ordered a large and it fits perfectly. There's about 10% off")
Recommended with %99.85

```
In [57]: 1 predict_recommendation("I don't like writing negative reviews but this one pissed me off the second I p
```

Not Recommended with %2.33

```
In [58]: 1 predict_recommendation("The cheapest material I've ever seen. It was like someone wove paper napkins fr
```

Not Recommended with %16.9

```
In [59]: 1 # Ref. 7 from now on  
2 predict_recommendation("I was so excited to receive this dress in the mail! The first day I wore it, I
```

Recommended with %91.75

```
In [60]: 1 predict_recommendation("The dress does not look like the dress pictures. The material seems cheaper and
```

Not Recommended with %2.57

```
In [61]: 1 predict_recommendation("I love this item it's was not dark blue like the picture but i love i it's very
```

Recommended with %73.58

```
In [62]: 1 predict_recommendation("This kaftan is NOT a silky material at all, it is a slightly transparent and du
```

Not Recommended with %4.94

```
In [ ]:
```

```
1
```

