# Text Classification - NLP Report

**Shubham Idekar**
College of Engineering
Northeastern University
Boston, MA
*idekar.s@northeastern.edu*

## Abstract

In this report, I explored the application of Multinomial Naive Bayes (NB) and Convolutional Neural Network (CNN) models for text classification in e-commerce, utilizing a robust dataset. Data preprocessing involved contraction expansion, punctuation removal, tokenization, lowercase conversion, elimination of numerical words, removal of stopwords, and stemming or lemmatization. The study aimed to categorize product descriptions into predefined classes. Multinomial NB served as a probabilistic baseline, while CNN, known for its capacity to capture intricate patterns in text, provided a more complex alternative.

## 1 Dataset

I used eCommerce dataset.



50425 rows × 2 columns

## 2 Preprocessing Data

### 2.1 Expanding Contraction

Contracted words are a common feature of natural language, especially in informal settings such as social media or messaging platforms.Contractions are shortened versions of words or phrases that are formed by combining two words and replacing one or more letters with an apostrophe. Examples of contractions include:

"can't" (from "cannot")
"won't" (from "will not")

```
df['Contractions'] = df['Text'].apply(lambda x: [contractions.fix(word) for word in x.split()])
df['No_contractions'] = [' '.join(map(str, l)) for l in df['Contractions']]
df.drop('Contractions',axis=1,inplace=True)
df.head()
```

| | Label | Text | No_contractions |
|---|---|---|---|
| 0 | Household | Paper Plane Design Framed Wall Hanging Motivat... | Paper Plane Design Framed Wall Hanging Motivat... |
| 1 | Household | SAF 'Floral' Framed Painting (Wood, 30 inch x ... | SAF 'Floral' Framed Painting (Wood, 30 inch x ... |
| 2 | Household | SAF 'UV Textured Modern Art Print Framed' Pain... | SAF 'UV Textured Modern Art Print Framed' Pain... |
| 3 | Household | SAF Flower Print Framed Painting (Synthetic, 1... | SAF Flower Print Framed Painting (Synthetic, 1... |
| 4 | Household | Incredible Gifts India Wooden Happy Birthday U... | Incredible Gifts India Wooden Happy Birthday U... |

Fig2.1: expanding contractions

## 2.2    Remove punctuations

Punctuation is often removed to simplify the analysis, and reduce the vocabulary size while preserving the meaningful content of the text.

We will use the punctuation library from the String package.

```
punc = string.punctuation
df['No_punc'] = df['No_contractions'].apply(lambda x: re.sub('[%s]' % re.escape(string.punctuation), '' , x))
df.head()
```

| | Label | Text | No_contractions | No_punc |
|---|---|---|---|---|
| 0 | Household | Paper Plane Design Framed Wall Hanging Motivat... | Paper Plane Design Framed Wall Hanging Motivat... | Paper Plane Design Framed Wall Hanging Motivat... |
| 1 | Household | SAF 'Floral' Framed Painting (Wood, 30 inch x ... | SAF 'Floral' Framed Painting (Wood, 30 inch x ... | SAF Floral Framed Painting Wood 30 inch x 10 i... |
| 2 | Household | SAF 'UV Textured Modern Art Print Framed' Pain... | SAF 'UV Textured Modern Art Print Framed' Pain... | SAF UV Textured Modern Art Print Framed Painti... |
| 3 | Household | SAF Flower Print Framed Painting (Synthetic, 1... | SAF Flower Print Framed Painting (Synthetic, 1... | SAF Flower Print Framed Painting Synthetic 135... |
| 4 | Household | Incredible Gifts India Wooden Happy Birthday U... | Incredible Gifts India Wooden Happy Birthday U... | Incredible Gifts India Wooden Happy Birthday U... |

Fig2.2 remove punctuations.

## 2.3    Tokenization

Tokenization is the process of breaking down text into individual words, phrases, or other meaningful elements, called tokens. We will use NLTK.word_tokenize() function to create a new column named "tokenized".

```
nltk.download('punkt')
df['Tokenized'] = df['No_punc'].apply(word_tokenize)
df.head()

[nltk_data] Downloading package punkt to C:\Users\Shubham
[nltk_data]    Idekar\AppData\Roaming\nltk_data...
[nltk_data]    Package punkt is already up-to-date!
```

| | Label | Text | No_contractions | No_punc | Tokenized |
|---|---|---|---|---|---|
| 0 | Household | Paper Plane Design Framed Wall Hanging Motivat... | Paper Plane Design Framed Wall Hanging Motivat... | Paper Plane Design Framed Wall Hanging Motivat... | [Paper, Plane, Design, Framed, Wall, Hanging, ... |
| 1 | Household | SAF 'Floral' Framed Painting (Wood, 30 inch x ... | SAF 'Floral' Framed Painting (Wood, 30 inch x ... | SAF Floral Framed Painting Wood 30 inch x 10 i... | [SAF, Floral, Framed, Painting, Wood, 30, inch... |
| 2 | Household | SAF 'UV Textured Modern Art Print Framed' Pain... | SAF 'UV Textured Modern Art Print Framed' Pain... | SAF UV Textured Modern Art Print Framed Painti... | [SAF, UV, Textured, Modern, Art, Print, Framed... |
| 3 | Household | SAF Flower Print Framed Painting (Synthetic, 1... | SAF Flower Print Framed Painting (Synthetic, 1... | SAF Flower Print Framed Painting Synthetic 135... | [SAF, Flower, Print, Framed, Painting, Synthet... |
| 4 | Household | Incredible Gifts India Wooden Happy Birthday U... | Incredible Gifts India Wooden Happy Birthday U... | Incredible Gifts India Wooden Happy Birthday U... | [Incredible, Gifts, India, Wooden, Happy, Birt... |

Fig2.3 Tokenization

## 2.4    Convert to Lower Case

All the alphabetic characters in a text are transformed to their corresponding lower case representation to reduce the vocabulary size and avoid duplication of words during text analysis.

```
df['Lower'] = df['Tokenized'].apply(lambda x: [text.lower() for text in x])
df.head()
```

| | Label | Text | No_contractions | No_punc | Tokenized | Lower |
|---|---|---|---|---|---|---|
| 0 | Household | Paper Plane Design Framed Wall Hanging Motivat... | Paper Plane Design Framed Wall Hanging Motivat... | Paper Plane Design Framed Wall Hanging Motivat... | [Paper, Plane, Design, Framed, Wall, Hanging, ... | [paper, plane, design, framed, wall, hanging, ... |
| 1 | Household | SAF 'Floral' Framed Painting (Wood, 30 inch x ... | SAF 'Floral' Framed Painting (Wood, 30 inch x ... | SAF Floral Framed Painting Wood 30 inch x 10 i... | [SAF, Floral, Framed, Painting, Wood, 30, inch... | [saf, floral, framed, painting, wood, 30, inch... |
| 2 | Household | SAF 'UV Textured Modern Art Print Framed' Pain... | SAF 'UV Textured Modern Art Print Framed' Pain... | SAF UV Textured Modern Art Print Framed Painti... | [SAF, UV, Textured, Modern, Art, Print, Framed... | [saf, uv, textured, modern, art, print, framed... |
| 3 | Household | SAF Flower Print Framed Painting (Synthetic, 1... | SAF Flower Print Framed Painting (Synthetic, 1... | SAF Flower Print Framed Painting Synthetic 135... | [SAF, Flower, Print, Framed, Painting, Synthet... | [saf, flower, print, framed, painting, synthet... |
| 4 | Household | Incredible Gifts India Wooden Happy Birthday U... | Incredible Gifts India Wooden Happy Birthday U... | Incredible Gifts India Wooden Happy Birthday U... | [Incredible, Gifts, India, Wooden, Happy, Birt... | [incredible, gifts, india, wooden, happy, birt... |

57

## *2.5      Remove words containing digits*

59   Eliminating words that contain numeric characters from text analysis to reduce noise and improve
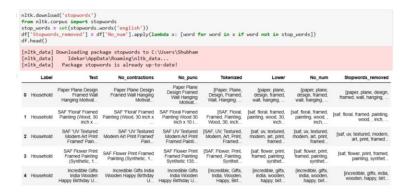60   the accuracy of language models. We will eliminate these words using Regular Expression.



61

## *2.6      Remove stopwords*

63   Process of eliminating common words such as "the", "a", "an", and "in" from text to reduce the
64   dimensionality of the data, and to focus on the more meaningful words that carry the essence of
65   the text.

66   We will use the stopwords library from the nltk module.



67
68                                    Fig 2.6 Removed stopwords

69

## *2.7      Stemming or Lemmatization*

71   Stemming and lemmatization are two techniques used in NLP to normalize words by reducing them
72   to their base or root form; stemming chops off the end of words, while lemmatization uses a
73   vocabulary and morphological analysis to reduce words to their canonical form.

74   Stemming: The stem of "running" is "run". Using a stemming algorithm, "running", "runs", and
75   "runner" would all be reduced to the stem "run".

76   Lemmatization: The lemma of "running" is "run". Using a lemmatization algorithm, "running" and
77   "runs" would be reduced to "run", while "runner" would be reduced to "run" as well, but only if the
78   context suggests that it is being used as a verb. We will apply parts of speech tags, in other words,
79   determine the part of speech (ie. noun, verb, adverb, etc.) for each word.

80

81   There are various stemmers and one lemmatizer in NLTK, the most common being:

82   Porter Stemmer from Porter (1980)

83   Wordnet Lemmatizer (port of the Morphy: https://wordnet.princeton.edu/man/morphy.7WN.html)

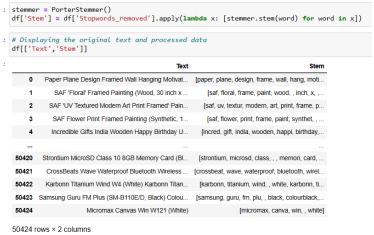84   Action : We will apply NLTK's Porter Stemmer within our trusty list comprehension.

```
: stemmer = PorterStemmer()
  df['Stem'] = df['Stopwords_removed'].apply(lambda x: [stemmer.stem(word) for word in x])

: # Displaying the original text and processed data
  df[['Text','Stem']]
```

| | Text | Stem |
|---|---|---|
| 0 | Paper Plane Design Framed Wall Hanging Motivat... | [paper, plane, design, frame, wall, hang, moti... |
| 1 | SAF 'Floral' Framed Painting (Wood, 30 inch x ... | [saf, floral, frame, paint, wood, , inch, x, ,... |
| 2 | SAF 'UV Textured Modern Art Print Framed' Pain... | [saf, uv, textur, modern, art, print, frame, p... |
| 3 | SAF Flower Print Framed Painting (Synthetic, 1... | [saf, flower, print, frame, paint, synthet, , ... |
| 4 | Incredible Gifts India Wooden Happy Birthday U... | [incred, gift, india, wooden, happi, birthday,... |
| ... | ... | ... |
| 50420 | Strontium MicroSD Class 10 8GB Memory Card (Bl... | [strontium, microsd, class, , , memori, card, ... |
| 50421 | CrossBeats Wave Waterproof Bluetooth Wireless ... | [crossbeat, wave, waterproof, bluetooth, wirel... |
| 50422 | Karbonn Titanium Wind W4 (White) Karbonn Titan... | [karbonn, titanium, wind, , white, karbonn, ti... |
| 50423 | Samsung Guru FM Plus (SM-B110E/D, Black) Colou... | [samsung, guru, fm, plu, , black, colourblack,... |
| 50424 | Micromax Canvas Win W121 (White) | [micromax, canva, win, , white] |

50424 rows × 2 columns

Fig 2.7 Stemming or lemmatization

# 3        Baseline Model- Multinomial Naïve Beyes

Here's a brief overview of how Multinomial Naive Bayes works for text classification:

*Bag-of-Words Representation:*

The first step is to represent the text data as a "bag of words," disregarding the order of words but considering their frequency. Each document is represented as a vector of word counts, where the elements of the vector correspond to the frequency of each word in the document.

*Vocabulary Building:*

The algorithm builds a vocabulary from the entire dataset, containing all unique words present in the corpus.

*Probability Estimation:*

MNB estimates the probability of each word occurring in a document for each class. It calculates the likelihood of observing each word given the class.

*Prior Probability:*

MNB also calculates the prior probability of each class, representing the likelihood of a document belonging to a particular class without considering the words.

*Naive Bayes Assumption:*

The "naive" assumption in Naive Bayes is that the features (words) are conditionally independent given the class label. While this assumption simplifies the model, it may not always hold true in practice.

*Classification Decision:*

Using Bayes' theorem, the algorithm calculates the posterior probability of each class given the document. The class with the highest posterior probability is assigned as the predicted class for the document.

**Baseline Model - Naive Bayes**

- Text Vectorization
  In order to perform machine learning on text data, we must transform the documents into vector representations. In natural language processing, **text vectorization** is the process of converting words, sentences, or even larger units of text data to numerical vectors.

```python
In [21]: from sklearn.model_selection import train_test_split

# Split the dataset into a training set and a testing set while preserving the class distribution
X = df['Stem']  # Features
y = df['Label']  # Target variable

# Use stratified sampling to ensure the same class distribution in both sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y, random_state=42)
```

```python
In [22]: from sklearn.feature_extraction.text import TfidfVectorizer

# Convert the lists of words into strings
X_train = X_train.apply(lambda word_list: ' '.join(word_list))
X_test = X_test.apply(lambda word_list: ' '.join(word_list))

tfidf_vectorizer = TfidfVectorizer()
X_train = tfidf_vectorizer.fit_transform(X_train)
X_test = tfidf_vectorizer.transform(X_test)
```

```python
In [23]: from sklearn.naive_bayes import MultinomialNB

nb_classifier = MultinomialNB()
nb_classifier.fit(X_train, y_train)
```

```
Out[23]: ▾ MultinomialNB
         MultinomialNB()
```

Fig3.1: Implementing of Multinomial NB

*Accuracy* and ***Prediction*** of this model :



```python
In [25]: from sklearn.metrics import classification_report, accuracy_score

accuracy = accuracy_score(y_test, predictions)
report = classification_report(y_test, predictions)
print(f"Accuracy: {accuracy}")
print(report)
```

```
Accuracy: 0.9411998016856717
                        precision    recall  f1-score   support

                 Books       0.98      0.92      0.95      2364
 Clothing & Accessories       0.98      0.94      0.96      1734
           Electronics       0.96      0.90      0.93      2124
             Household       0.90      0.98      0.94      3863

              accuracy                           0.94     10085
             macro avg       0.95      0.93      0.94     10085
          weighted avg       0.94      0.94      0.94     10085
```

Fig3.2: Model Performance by metrics.

Multinomial Naive Bayes Model Evaluation:

Precision:
Books: 0.98 - Among the instances predicted as "Books," 98% are correctly classified.
Clothing & Accessories: 0.98 - 98% of instances predicted as "Clothing & Accessories" are correct.
Electronics: 0.96 - 96% of instances predicted as "Electronics" are correct.
Household: 0.90 - 90% of instances predicted as "Household" are correct.

Recall:
Books: 0.92 - The model correctly identifies 92% of the actual instances of "Books."
Clothing & Accessories: 0.94 - 94% of instances of "Clothing & Accessories" are correctly identified.
Electronics: 0.90 - The model captures 90% of instances of "Electronics."
Household: 0.98 - An impressive 98% of instances of "Household" are correctly identified.

F1-Score:
Books: 0.95 - The harmonic mean of precision and recall for "Books."
Clothing & Accessories: 0.96 - The F1-score for "Clothing & Accessories."
Electronics: 0.93 - The harmonic mean of precision and recall for "Electronics."
Household: 0.94 - The F1-score for "Household."

Support:
Books: 2364 - There are 2364 instances of "Books" in the test set.
Clothing & Accessories: 1734 - There are 1734 instances of "Clothing & Accessories."
Electronics: 2124 - There are 2124 instances of "Electronics."

145 Household: 3863 - There are 3863 instances of "Household."
146
147 Accuracy:
148 The overall accuracy of the MultinomialNB model is approximately 94%, meaning it correctly
149 classifies instances around 94% of the time.
150
151 Macro Avg and Weighted Avg:
152 Macro Avg: 0.94 - The average precision, recall, and F1-score across all classes without
153 considering class imbalance.
154 Weighted Avg: 0.94 - Similar to macro avg, but considering the number of samples for each class,
155 giving more weight to classes with more instances.
156
157 Interpretation:
158 The MultinomialNB model shows good performance across all classes. It performs particularly
159 well in correctly identifying instances of "Books," "Clothing & Accessories," and "Household."
160 The weighted average considers the class imbalance, providing a balanced overview of model
161 performance.
162
163
## 4      Advanced model - CNN

166 Convolutional Neural Networks (CNNs) are primarily known for their success in
167 image-related tasks, but they can also be adapted for text classification. Here's a
168 simplified explanation of how a CNN model for text classification works:
169
170 *Input Representation*:
171 Text data is initially represented as numerical vectors, typically using techniques
172 like word embeddings (e.g., Word2Vec, GloVe). Each word is represented by a
173 vector, and a sequence of these vectors is used as the input.
174
175 *Convolutional Layers*:
176 Convolutional operations, which are highly effective in capturing local patterns,
177 are applied to the input sequence. This involves using filters or kernels of fixed
178 size that slide over the input.
179 The convolutional layer detects specific patterns or features in the input. In text,
180 these features can represent n-grams or local word patterns.
181
182 *Activation and Pooling*:
183 After the convolution, an activation function (e.g., ReLU) is applied element-wise
184 to introduce non-linearity.
185 Pooling layers (often max pooling) follow, reducing the dimensionality of the
186 features while retaining the most salient information. This helps the model focus
187 on the most important features.
188
189 *Flattening and Fully Connected Layers*:
190 The output from the convolutional and pooling layers is flattened into a one-
191 dimensional vector.
192 Fully connected layers are then used to combine the learned features and make
193 predictions. These layers can capture global relationships in the data.
194
195 *Output Layer*:
196 The final layer consists of one or more neurons representing the classes in the
197 classification task. A softmax activation function is often used to convert the
198 network's raw output into class probabilities.
199 *Training*:
200 The model is trained using labeled data with a specified loss function (e.g.,
201 categorical cross-entropy). The weights of the network are updated through
202 backpropagation and optimization algorithms (e.g., Adam, SGD).

203
204 *Prediction*:
205 Once trained, the model can predict the class of a new text by passing it through
206 the trained network, and the class with the highest probability is considered the
207 predicted class.

208
```python
In [24]: # Split the data into training and testing sets
train_data, test_data = train_test_split(df, test_size=0.2, random_state=42, stratify=df['Label'])

# Advanced Model: Convolutional Neural Network (CNN)
label_encoder = LabelEncoder()
y_train_encoded = label_encoder.fit_transform(train_data['Label'])
y_test_encoded = label_encoder.transform(test_data['Label'])

tokenizer = Tokenizer()
tokenizer.fit_on_texts(train_data['Stem'])
X_train_sequences = tokenizer.texts_to_sequences(train_data['Stem'])
X_test_sequences = tokenizer.texts_to_sequences(test_data['Stem'])

max_sequence_length = max(max(len(seq) for seq in X_train_sequences), max(len(seq) for seq in X_test_sequences))
X_train_padded = pad_sequences(X_train_sequences, maxlen=max_sequence_length, padding='post')
X_test_padded = pad_sequences(X_test_sequences, maxlen=max_sequence_length, padding='post')

embedding_dim = 50
vocab_size = len(tokenizer.word_index) + 1

cnn_model = Sequential()
cnn_model.add(Embedding(input_dim=vocab_size, output_dim=embedding_dim, input_length=max_sequence_length))
cnn_model.add(Conv1D(filters=128, kernel_size=5, activation='relu'))
cnn_model.add(GlobalMaxPooling1D())
cnn_model.add(Dense(64, activation='relu'))
cnn_model.add(Dropout(0.5))
cnn_model.add(Dense(len(label_encoder.classes_), activation='softmax'))

cnn_model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy'])
cnn_model.fit(X_train_padded, y_train_encoded, epochs=5, batch_size=64, validation_split=0.2)

cnn_accuracy = cnn_model.evaluate(X_test_padded, y_test_encoded, verbose=0)[1]
print("\nAdvanced Model (Convolutional Neural Network) Results:")
print(f"Accuracy: {cnn_accuracy}")


Epoch 1/5
505/505 [==============================] - 407s 802ms/step - loss: 0.4292 - accuracy: 0.8483 - val_loss: 0.1321 - val_accuracy: 0.9646
Epoch 2/5
505/505 [==============================] - 419s 829ms/step - loss: 0.0933 - accuracy: 0.9787 - val_loss: 0.0985 - val_accuracy: 0.9746
Epoch 3/5
505/505 [==============================] - 416s 824ms/step - loss: 0.0387 - accuracy: 0.9907 - val_loss: 0.1047 - val_accuracy: 0.9747
Epoch 4/5
505/505 [==============================] - 419s 829ms/step - loss: 0.0207 - accuracy: 0.9949 - val_loss: 0.1215 - val_accuracy: 0.9763
Epoch 5/5
505/505 [==============================] - 423s 837ms/step - loss: 0.0125 - accuracy: 0.9969 - val_loss: 0.1342 - val_accuracy: 0.9761

Advanced Model (Convolutional Neural Network) Results:
Accuracy: 0.9734258651733398
```

209

210 Evaluation metrics for this model:

```
316/316 [==============================] - 17s 54ms/step
Classification Report for CNN Model:
                        precision    recall  f1-score   support

                 Books       0.97      0.97      0.97      2364
 Clothing & Accessories       0.98      0.98      0.98      1734
           Electronics       0.97      0.97      0.97      2124
             Household       0.97      0.97      0.97      3863

              accuracy                           0.97     10085
             macro avg       0.97      0.97      0.97     10085
          weighted avg       0.97      0.97      0.97     10085
```

211

212 CNN Model Evaluation:

213 *Precision*:
214 Precision is the ratio of correctly predicted positive observations to the
215 total predicted positives. For each class, precision is calculated as TP / (TP
216 + FP), where TP is the number of true positives and FP is the number of
217 false positives. In your report, precision values range from 0.97 to 0.98,
218 indicating high precision for all classes.

219 *Recall*:
220 Recall (or sensitivity or true positive rate) is the ratio of correctly
221 predicted positive observations to the all observations in the actual class.
222 For each class, recall is calculated as TP / (TP + FN), where TP is the
223 number of true positives and FN is the number of false negatives. In your
224 report, recall values range from 0.97 to 0.98, indicating high recall for all
225 classes.

226 *F1-Score*:

227  F1-score is the weighted average of precision and recall. It considers both
228  false positives and false negatives. For each class, F1-score is calculated as
229  2 * (precision * recall) / (precision + recall). In your report, F1-score
230  values range from 0.97 to 0.98, indicating a good balance between precision
231  and recall for all classes.

232  *Support*:
233  Support is the number of actual occurrences of the class in the specified
234  dataset. It gives you an idea of how many samples in your test set belong to
235  each class.

236  *Overall Accuracy*:
237  The accuracy is the overall correctly predicted instances divided by the
238  total instances. In your case, the overall accuracy is around 97%, indicating
239  the proportion of correctly classified instances.

240  *Macro Avg and Weighted Avg*:
241  Macro avg is the average of precision, recall, and F1-score across all
242  classes without considering class imbalance. Weighted avg is the same as
243  macro avg, but it considers the number of samples for each class, giving
244  more weight to classes with more instances.

245  *Interpretation*:
246  The high precision, recall, and F1-score values for each class and the
247  overall accuracy indicate that our CNN model is performing well on the test
248  set.
249
250  **References**

251  [1]    https://www.analyticsvidhya.com/blog/2020/12/understanding-text-classification-in-nlp-with-
252  movie-review-example-example/

253  [2] https://www.kaggle.com/code/firozchowdury/data-pre-processing-e-commerce-dataset

254  [3] https://www.kaggle.com/code/bekkarmerwan/ecommerce-text-classification