

# A Survey on Usage of Vectorizers for Textual Data in Exploratory Data Analysis (EDA)-Generative AI

DEVARAM SHARATHCHNDRA

Software Engineer, Machine Learning Enthusiast

**Abstract**— In the early age of Artificial Intelligence (AI) and Machine Learning (ML) domain, mostly training of ML models are depended on numerical data to classify, predict or generate. In today’s world we achieved the state “Machine models” can interact with human in a pure form of humanized text. Natural Language Processing (NLP) is the growing domain where it interacts with human in a way of speech recognition, text classification and text generation. The present era is experiencing prompt-based AI, where we can generate new images with a simple text prompt input or can generate a professional video or chat bot types models for virtual assistance. Simultaneously we are interacting with speech with a machine. The core technology behind this textual input is vectorizing the text data. When we interact with ML model with a speech input, in the background-the speech is converted into a textual format and then vectorized for prediction or generation to produce output. Based on the produced output the output layer can interact with human according to the choice provided by the end user weather it is belonging to NLP or Text generation transformer type model. The best example for humanized text generation model we are experiencing in today’s technology era are Google’s Gemini and Open AI’s Generative Pre-Trained Transformer (GPT) model. Vectorizers are the main technology behind these text transformation and analyzation models. The main amin these vectorizers re to improve machine learning model accuracy and reducing computational complexity of a ML model. NLP use multilayered neural networks for a Deep Learning (DL) model. Before feeding the first input layer with this textual data, we are using this vectorizers concept while training the deep learning model. Vectorization concept is involved in feature extraction and these will include different type of vectorizers. In this survey paper we discussed most of the vectorizers in section wise. In the I. Introduction section, I am going to introduce the concepts of vectors and what are different types of vectors available to use for machine learning model. From the section II. Core Technology, I’ll explain how we use vectorizers for a Machine Learning, Deep Learning and Transformer models to train. From the final section III. Results, difference between all type of vectorizers are concluded.

**Index Terms**—Vectorizers, Machine Learning (ML), Deep Learning, NLP, Transformers, Artificial Intelligence (AI).

## I. INTRODUCTION

The term vectorizer itself conveys a primary meaning “vector” which means it has some directional things. We can define vectors with its close relative term Matrix. Similarly vector also has a size. As we are considering textual data for a machine learning model, we converting text to numbers with this vectorizers. The process of converting textual data to numerical data is called as “Vectorizing”. Fitting of a ML model for a textual data is achieved with the help of vectorizers. The textual data may include a paragraph or a group of words or some letters to categorize. The below image represents the concept of vectorization.

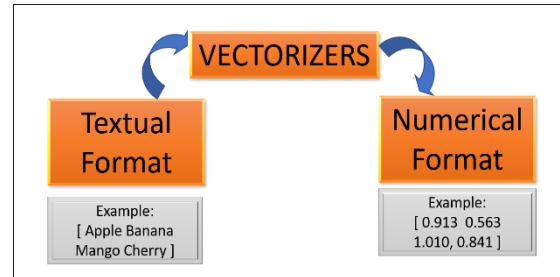


Fig 1. Concept of Vectorizers.

Types of Vectorizers:

Vectorizers<sup>[1-8]</sup> are primary method to convert textual data to numerical data. When we have a different type of data, we deal with different types of vectorizers.

1. Count Vectorizers:

The input textual data has a structure of sentences, we can process with Count Vectorizers. In general count vectorizers are able to differentiate vocabulary from the sentence input. Once vectorizer was initialized with count vectorizers module from sci-kit learn, we are fitting the sentence into vectorizers. From the below example we can clearly observe what is

happening with count vectorizers for a sentence structure.

Example sentence: “Hi hello how are you, I am from India. We are using vectorizers to categorize text to numerical.”

After vectorization process, we can observe a numerical array format for the specific sentence. The numerical array is size of One row and 15 columns matrix of class integer with NumPy. The size of the matrix is [1 x 15]. Numerical array is as this format [[2, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 1, 1, 1, 1]]. The logic behind these count vectorizers are, primarily these Count Vectorizers differentiate unique vocabulary from the input sentence and fixes the size of array. For our input sentence we have 15 unique words. So, the size of output NumPy array is 15. Once vocabulary is fixed, it will compare the frequency of the word with in sentence.

```

sentence = """Hi hello how are you,
            Iam from India. We are using vectorizers
            to cateorize text to numericals."""

vect_data = cv.fit_transform([sentence])

array(['are', 'cateorize', 'from', 'hello', 'hi', 'how', 'iam', 'india',
       'numericals', 'text', 'to', 'using', 'vectorizers', 'we', 'you'],
      dtype=object)

numerical_array
array([[2, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 1, 1, 1, 1]])
    
```

Fig 2. Representation of Count Vectorizers for a Single Sentence.

Count Vectorizers for a paragraph: While we are dealing with a paragraph count vectorizers collects all the vocabulary in all sentences and makes a numerical array with vocabulary length. This will create a multi-dimensional NumPy array. The example of a paragraph is explained. Example paragraph is “I love India”. “I am a python programmer”. “Programming is a fun task”. We got the numerical array as [3 x9] matrix. [[0 0 1 0 1 0 0 0 0][1 0 0 0 0 1 0 1 0][0 1 0 1 0 0 1 0 1]].

```

[37] documents = [
        "I love India",
        "I am a python programmer",
        "Programming is a fun task"
    ]

Document-Term Matrix:
[[0 0 1 0 1 0 0 0 0]
 [1 0 0 0 0 1 0 1 0]
 [0 1 0 1 0 0 1 0 1]]
    
```

Fig 3. Usage of count vectorizers for a paragraph.

## 2. Tokenizers

The world of textual data which is dealing with Natural Language Processing (NLP) uses tokenizers concept. The tokenizers are very similar to count vectorizers. Tokenizers are available in both Natural Language Toolkit (NLTK) and Keras module. From nltk, tokenizers first convert paragraphs into sentences and then sentences to words to numerical. Tokenizers from keras module arranges all vocabulary in an order through frequency count of a word in total vocabulary. The highest frequency word gets lowest number and vice versa. Here is the example usage of a paragraph “Hi all, how are you. We are discussing about the topic tokenizers. We are considering both the module nltk and keras preprocessing module.” The result is paragraph is divided into sentences and then words with the help of word tokenizer from nltk. The result is ['Hi all, how are you.', 'We are discussing about the topic tokenizers.', 'We are considering both the module nltk and keras preprocessing module.'], after the word tokenizers we need to consider for further evaluation. The below picture shows the details.

```

sent_tokenize(paragraph)

['Hi all, how are you.',
 'We are discussing about the topic tokenizers.',
 'We are considering both the module nltk and keras preprocessing module.']
    
```

Fig 4. Sentences defined for toekniztion.

The word tokenizer result is ['Hi', 'all', ',', 'how', 'are', 'you.', 'We', 'are', 'discussing', 'about', 'the', 'topic', 'tokenizers.', 'We', 'are', 'considering', 'both', 'the', 'module', 'nltk', 'and', 'keras', 'preprocessing',





```

# Define the model
model = Sequential()
model.add(Embedding(input_dim=len(tokenizer.word_index) + 1,
                    output_dim=128, input_length=max_length))
model.add(Flatten())
model.add(Dense(units=128, activation="relu"))
model.add(Dense(units=len(one_hot_labels[0]), activation="softmax"))

model.compile(optimizer="adam", loss="categorical_crossentropy", metrics=["accuracy"])
model.fit(x=train, y=train, epochs=10, batch_size=32, validation_data=(x_test, y_test))
    
```

```

Epoch 1/10
1/1 ----- 4s 4s/step - accuracy: 0.5000 - loss: 1.3815 - val_accuracy: 0.5000 - val_loss: 1.3936
Epoch 2/10
1/1 ----- 0s 353ms/step - accuracy: 0.8333 - loss: 1.1709 - val_accuracy: 0.0000e+00 - val_loss: 1.4010
Epoch 3/10
1/1 ----- 0s 96ms/step - accuracy: 0.8333 - loss: 1.0107 - val_accuracy: 0.0000e+00 - val_loss: 1.4050
Epoch 4/10
1/1 ----- 0s 66ms/step - accuracy: 0.8333 - loss: 0.8638 - val_accuracy: 0.0000e+00 - val_loss: 1.4009
Epoch 5/10
1/1 ----- 0s 140ms/step - accuracy: 0.8333 - loss: 0.7294 - val_accuracy: 0.0000e+00 - val_loss: 1.4091
Epoch 6/10
1/1 ----- 0s 65ms/step - accuracy: 0.8333 - loss: 0.6077 - val_accuracy: 0.0000e+00 - val_loss: 1.4041
Epoch 7/10
1/1 ----- 0s 67ms/step - accuracy: 0.8333 - loss: 0.4938 - val_accuracy: 0.0000e+00 - val_loss: 1.3911
Epoch 8/10
1/1 ----- 0s 134ms/step - accuracy: 0.8333 - loss: 0.4056 - val_accuracy: 0.0000e+00 - val_loss: 1.3679
Epoch 9/10
1/1 ----- 0s 64ms/step - accuracy: 1.0000 - loss: 0.3248 - val_accuracy: 0.0000e+00 - val_loss: 1.3341
Epoch 10/10
1/1 ----- 0s 66ms/step - accuracy: 1.0000 - loss: 0.2555 - val_accuracy: 0.0000e+00 - val_loss: 1.2915
<keras.src.callbacks.history.History at 0x797b1d64a60b>
    
```

Fig 13. Deep Learning Model with Tokenizer Concept and Model Fitting.

### 3. Transformer Model with Tokenizers concept.

Transformer models are extended version of deep learning, which are able to generate new textual data based on the trained data. When new text prompt is given to this model, it processes the data and produces the text generated from trained data. A simple transformer model is considered below for this survey analysis. Consider a simple ‘BERT’ transformer model with tokenizer. BERT is a Bidirectional Representations from Transformers. This model uses tokenizers from hugging face platform.

```

[23] sample_text = "This is a sample text for survey paper for vectorizers and tokenizers"
[24] tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
tokenizer_config.json: 100% |#####| 48.0/48.0 [00:00<00:00, 1.18kBi/s]
vocab.txt: 100% |#####| 232k/232k [00:00<00:00, 3.93MB/s]
tokenizer.json: 100% |#####| 466k/466k [00:00<00:00, 4.09MB/s]
config.json: 100% |#####| 570/570 [00:00<00:00, 7.76kBi/s]
input_layer = tokenizer([sample_text, return_tensors="pt"])
    
```

Fig 14. Transformer model from Hugging face platforms.

### III. RESULTS

As we discussed different types of vectorizers and their application across models, choosing the exact vectorizer or tokenizers purely depends on the use case of a machine learning model. If we are about to generate series of sentences, tokenizers would be a good choice in some cases. For word-based gaming models word vectorizers or count vectorizers are good

choice. Below image represent comparison of count vectorizer and tokenizers.

First image represents usage of count vectorizers and the second image shows usage of tokenizers in for the same text documents.

```

[31] # Sample text data
documents = [
    "I love my country",
    "India is my country",
    "India has highest population among all countries",
    "I love cricket",
]
    
```

Fig 15. Sample text documents for survey analysis.

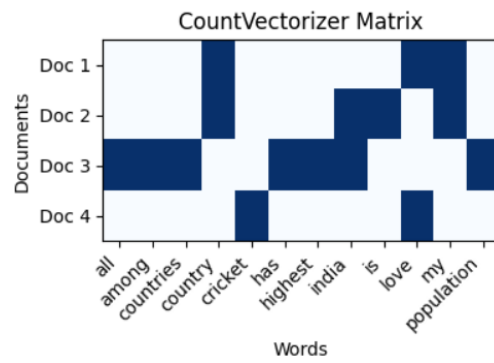


Fig 16. Count vectorizer plots for text for four documents

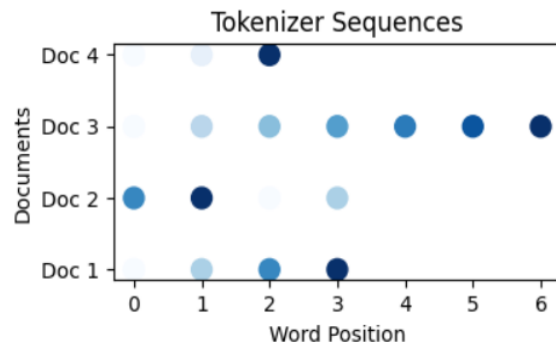


Fig 17. Tokenizers from Keras module in for same four text documents.

### CONCLUSION

This survey paper discussed about effective usage of vectorizers from small scale ML models to large language models (Transformers). Different types of text to numeric conversation helpful tin training machine learning models.

REFERENCES

- [1] Machine Learning in python via scikit learn official from <https://scikit-learn.org/stable/>
- [2] Transformer models from hugging face platform, <https://huggingface.co/>
- [3] Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. *Advances in Neural Information Processing Systems*, pp. 1–9
- [4] Harris, C.R., Millman, K.J., van der Walt, S.J. et al. Array programming with NumPy. *Nature* 585, 357–362 (2020). DOI: 10.1038/s41586-020-2649-2. (Publisher link).
- [5] The pandas development team, ‘pandas-dev/pandas: Pandas’. Zenodo, Apr. 10, 2024. doi: 10.5281/zenodo.10957263.
- [6] Buitinck, L. et al., 2013. API design for machine learning software: experiences from the scikit-learn project. In *ECML PKDD Workshop: Languages for Data Mining and Machine Learning*. pp. 108–122.
- [7] Chollet F, others. Keras [Internet]. GitHub; 2015. Available from: <https://github.com/fchollet/keras>
- [8] Nazzal, Fayez & Abualrob, Ahmad & Husni, Ahmad. (2022). Classification of Item's category using the CountVectorizer, The Linear Support Vector Machine Algorithm and Python's sklearn module. 10.13140/RG.2.2.24423.32169.
- [9] Yang, Jinbiao. (2024). Rethinking tokenization: Crafting better tokenizers for large language models. *International Journal of Chinese Linguistics*. 11. 94-109. 10.1075/ijchl.00023.yan.
- [10] (2024). Large language models (LLMs): survey, technical frameworks, and future challenges. *Artificial Intelligence Review*. 57. 10.1007/s10462-024-10888-y.
- [11] Yang, Jingfeng & Jin, Hongye & Tang, Ruixiang & Han, Xiaotian & Feng, Qizhang & Jiang, Haoming & Zhong, Shaochen & Yin, Bing & Hu, Xia. (2024). Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond. *ACM Transactions on Knowledge Discovery from Data*. 18. 10.1145/3649506.
- [12] Ananthaswamy, Anil. (2023). In AI, is bigger always better?. *Nature*. 615. 202-205. 10.1038/d41586-023-00641-w.
- [13] Raza, Muhammad & Meghji, Areej & Mahoto, Naeem & Al Reshan, Mana & Abosaq, Hamad & Sulaiman, Adel & Shaikh, Asadullah. (2024). Reading Between the Lines: Machine Learning Ensemble and Deep Learning for Implied Threat Detection in Textual Data. *International Journal of Computational Intelligence Systems*. 17. 10.1007/s44196-024-00580-y.
- [14] Ulah, Arif & Khan, Sundas & Mohd Nawi, Nazri. (2022). Review on sentiment analysis for text classification. *Multimedia Tools and Applications*. 82. 10.1007/s11042-022-14112-3.
- [15] Zhang, Tianyi & Ladhak, Faisal & Durmus, Esin & Liang, Percy & McKeown, Kathleen & Hashimoto, Tatsunori. (2024). Benchmarking Large Language Models for News Summarization. *Transactions of the Association for Computational Linguistics*. 12. 39-57. 10.1162/tacl\_a\_00632.