



ISSN: 0976-3376

Available Online at <http://www.journalajst.com>

ASIAN JOURNAL OF  
SCIENCE AND TECHNOLOGY

Asian Journal of Science and Technology  
Vol. 14, Issue, 06, pp. 12555-12560, June, 2023

## RESEARCH ARTICLE

# LABELING OF TEXT DATA USING AUTOENCODERS

\*Praveen Thenraj Gunasekaran, Selvakuberan Karuppasamy and Subhashini Lakshminarayanan

Artificial Intelligence, Accenture, Chennai, India

### ARTICLE INFO

#### Article History:

Received 03<sup>rd</sup> March, 2023  
Received in revised form  
26<sup>th</sup> April, 2023  
Accepted 14<sup>th</sup> May, 2023  
Published online 20<sup>th</sup> June, 2023

#### Keywords:

Unlabeled text data, Auto Labeling,  
AutoEncoders, Clustering.

### ABSTRACT

Machine learning has come a long way in solving business use cases that has remained a nightmare to human. Today machines learn data in ways like human, machine learning has matured so much that all it requires is data and it can solve any problem if the correct data is provided. Among the different learning techniques, we have in current ML world, supervised learning is a popular technique where the model learns from labeled dataset. The model tries to learn the pattern from the data and tries to correlate the independent and the dependent variable. But the challenge in real time is we don't have the readily available labeled data which applies to unstructured text as well. Given the volume of the text data available and the multiple sources available, it would take humongous efforts to label these text data manually. This has led to the rise of many unsupervised techniques to learn the data for solving use cases. However, in spite of numerous improvements in the domain of unsupervised learning, the supervised learning continues to be one of the preferred techniques for humans to train machines. The objective of this paper is to use AutoEncoders combined with clustering technique to label the unlabeled text training data when the number of classes for the dataset is known.

**Citation:** Praveen Thenraj Gunasekaran, Selvakuberan Karuppasamy and Subhashini Lakshminarayanan. 2023. "Labeling of text data using autoencoders", *Asian Journal of Science and Technology*, 14, (06), 12555-12560.

Copyright©2023, Praveen Thenraj Gunasekaran et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

## INTRODUCTION

In the era of deep learning, labelling the training data manually is a very tedious task, given the volume of training data that is being used. With the advent of machine learning which can solve many use cases in many domains, there must be techniques to solve its own problem of getting the data labeled for training purposes. There are also readily available labeling tools which can help labeling unlabeled dataset but still the reliability of these packages remains a question when it comes to critical business scenarios. In this paper, we propose a simple solution, based on autoencoder and clustering to solve the problem of labeling unlabeled text data. The solution consists of four parts 1. Embed the training dataset 2. Extract important features of the training dataset 3. Clustering of the lower dimension representation 4. Keyword identification from each cluster.

## BACKGROUND

**Text data labeling:** Text data is a form of unstructured data. There are various sources of text data especially with the advent of internet and social media, the volume of unstructured data available also has increased linearly. Increase in volume also means annotating this huge volume of text data involves huge amount of human effort. Since we are dealing with big data, human intervention for such a huge volume of data would result in more resource necessity, accuracy in annotation as different people with different perceptions would be involved and also increased cost. In the process of

continuous improvement, there has been some cool techniques semi-supervised learning that has been identified to solve the problem of unlabeled dataset. Unsupervised techniques also can be used to label the training data whereas semi-supervised techniques make use of a considerable portion of the training data that has already been labeled and uses them to learn and label the remaining dataset.

## LITERATURE REVIEW

In [1], the authors have used an autoencoder and clustering based technique to solve the problem of labeling image dataset. The authors have used MNIST dataset for this experiment. [7] A Siamese network-based architecture to derive the sentence embeddings of a given pair of sentences. This approach is a modified version of the pretrained BERT model, and it generates more relevant embeddings with much reduction in computation time as well. [4] uses Deep Autoencoders along with SVM as a classification layer for classifying the images. The authors have used MNIST dataset for this work and have obtained 99.8% accuracy. [8] This paper marked a new era in the domain of NLP. The authors realized the need for understanding the contextuality of the tokens in a sentence and came up with two architectures namely CBOW and Skipgram to generate word embeddings for English language that can be used across any tasks. [6] The authors in the paper have used K-Means algorithm as clustering technique for clustering the similar national anthems of different countries of the world. The authors have used TF-IDF as mechanism to extract the features from the documents and then used K-Means algorithm to cluster the documents. In [2], the authors have

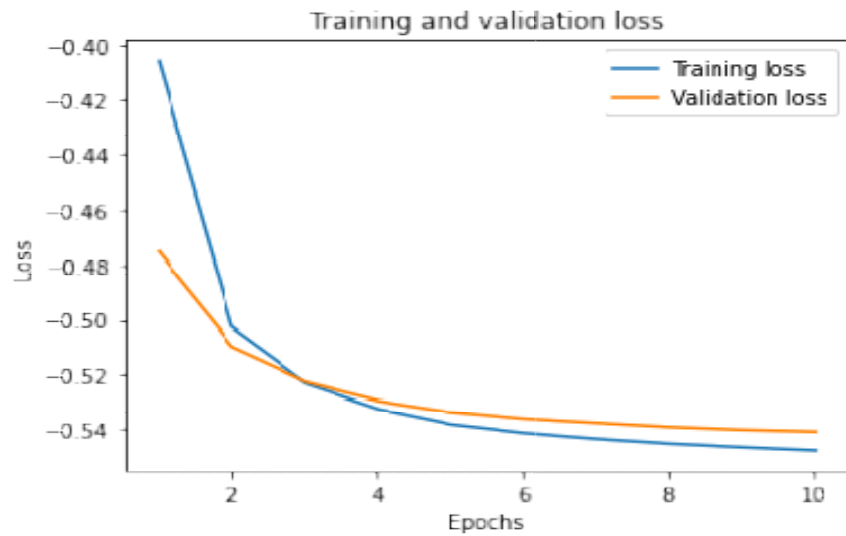
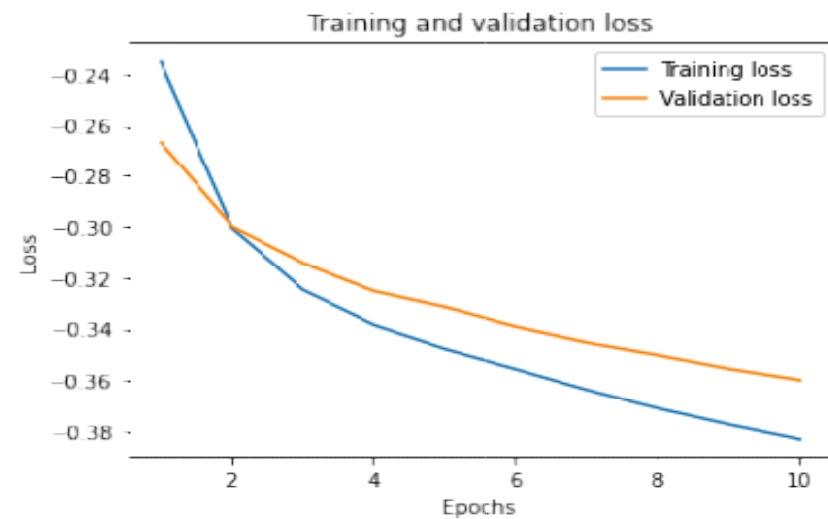






**Table 1. Comparison of metrics during different iterations of the experiment**

Number of Hidden Layers	Bottleneck Layer Dimension	Epochs	Class												Accuracy
			Business			Science/Technology			Sports			World			
			Precision	Recall	F1-Score	Precision	Recall	F1-Score	Precision	Recall	F1-Score	Precision	Recall	F1-Score	
1	128	10	0.4	0.59	0.47	0.18	0.15	0.16	0.26	0.25	0.26	0.01	0.01	0.01	0.25
	150	10	0.37	0.55	0.45	0.95	0.69	0.8	0.29	0.24	0.26	0.33	0.31	0.32	0.45
	175	10	0.23	0.2	0.21	0.95	0.67	0.79	0.35	0.52	0.42	0.34	0.31	0.33	0.43
6	16	10	0.03	0.03	0.03	0	0	0	0.46	0.74	0.57	0.03	0.02	0.03	0.2
7	150	10	0.88	0.75	0.81	0.88	0.9	0.89	0.32	0.34	0.33	0.37	0.39	0.38	0.6
9	100	10	0.01	0.01	0.01	0.94	0.68	0.79	0.39	0.61	0.47	0.07	0.05	0.06	0.34
9	150	10	0.05	0.05	0.05	0.01	0.01	0.01	0.37	0.45	0.41	0.02	0.01	0.01	0.13
21	10	10	0.09	0.25	0.14	0.21	0.06	0.09	0	0	0	0.2	0.2	0.2	0.13

**Fig. 11. Bottleneck layer dimension – 150, Number of Hidden Layers – 7****Fig. 12. Bottleneck layer dimension – 10, Number of Hidden Layers – 21**

and cluster them to create labels. We have run experiments with different iterations by varying the neuron size and hidden layer size and keeping the number of epochs constant. Our initial objective was to understand how efficient we can use this technique for labeling an unlabeled text dataset and we have achieved 0.6 accuracy using it. The primary scope of future work is to try this architecture for a domain specific dataset (eg. Banking, Insurance) to understand how well domain knowledge can be captured using this technique, as there will be limitations like limited labeled data for domain specific data. We would also like to extend the scope of this work to try different embedding techniques to generate the input for the autoencoder.

## REFERENCES

- [1] Paraskevi Nousi, Anastasios Tefas, "Self-supervised autoencoders for clustering and classification", Springer-Verlag GmbH Germany, part of Springer Nature 2018
- [2] Soodeh Hosseini, Zahra Asghari Varzaneh, "Deep text clustering using stacked AutoEncoder", Springer Science+Business Media, LLC, part of Springer Nature 2022
- [3] Dor Bank, Noam Koenigstein, Raja Giryes, "Autoencoders," arXiv:2003.05991v2 [cs.LG] 3 Apr 2021
- [4] Munmi Gogoi, Shahin Ara Begum, "Image Classification using Deep Autoencoders", 2017 IEEE International Conference on Computational Intelligence and Computing Research (ICCIC)
- [5] Zhenyu Yang, Xue Pang, "Research on Text Classification of Denoising Autoencoder Based on Additional Momentum and Adaptive Learning Rate", 2018 11th International Symposium on Computational Intelligence and Design (ISCID)
- [6] Chaman Lal et al., "Text Clustering using K-MEAN " , International Journal of Advanced Trends in Computer Science and Engineering, 10(4), July – August 2021, 2892 – 2897 2892
- [7] Nils Reimers, Iryna Gurevych, "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks", EMNLP 2019
- [8] Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean, "Efficient Estimation of Word Representations in Vector Space", arXiv:1301.3781v3 [cs.CL] 7 Sep 2013
- [9] Piotr Bojanowski, Edouard Grave, Armand Joulin, Tomas Mikolov, "Enriching Word Vectors with Subword Information", arXiv:1607.04606v2 [cs.CL] 19 Jun 2017

\*\*\*\*\*