

Performance Evaluation of Dimensionality Reduction Techniques on Big Data using Machine Learning Algorithms

Venkateswara Gupta Pola¹, Pranahith Babu Yarra², Ramya. A³, Sai Suraj Karra⁴

¹Venkateswara Gupta Pola

Dept. of Computer Science Engineering, SASTRA University

²Pranahith Babu Yarra

Dept. of Mechanical Engineering (Mechatronics), Mahatma Gandhi Institute of Technology

³Ramya. A

Dept. of Computer Science Engineering, Marri Laxman Reddy Institute of Technology and Management

⁴Sai Suraj Karra

Dept. of Mechanical Engineering (Mechatronics), Mahatma Gandhi Institute of Technology

Abstract - Large volumes of data are generated due to digitization which in-turn created huge datasets. To find specific patterns hidden in the attributes of these datasets, several ML algorithms are used. But in these large datasets, certain attributes may not be helpful or may not affect the prediction result. Ignoring irrelevant attributes reduces burden on ML algorithms and decreases cost and computation power which can be done using dimensionality reduction techniques.

In this project two prominent Dimensionality Reduction algorithms are used namely Linear Discriminant Analysis (LDA) and Principal Component Analysis (PCA) on Intrusion detection dataset with ML Techniques like DT, NB, SVM and RF classifiers. The project was completed with:

1. 4 Machine learning algorithms namely DT, SVM, NB, RF with two Dimensionality reduction algorithms PCA and LDA respectively.
2. PySpark which is a Big Data tool used in this project for implementing the machine learning and dimensionality reduction algorithms.
3. Encoders like Label Encoder and One Hot Encoder for converting categorical data to binary data making it suitable for ML algorithms.

Key Words: Machine learning, PCA, LDA, Algorithms.

1. INTRODUCTION

In the past 15 years, Machine Learning has been one among the fastest growing technologies. It has various applications in different fields, such as computer vision, computational biology, data science, healthcare, banking, criminal identification, market prediction, etc. ML enables a computing machine to learn from a huge sample of information and predict the trends within the raw data. In

various research areas, ML algorithms are used to predict hidden patterns and identify the test data to generate accurate results. The scope of ML classification models in the area of cybersecurity is constantly growing. They have proven to be very effective in detecting numerous threats and malicious content.

1.1 WHAT IS INTRUSION DETECTION?

We are being reliant day by day on networks and computer technologies. This raises the need for stable networks. For data privacy, confidentiality and availability, we have to strengthen computer network protection. But these do not stop the detection of intrusion. It is important to protect insecure computing systems and networks only to avoid the possibility of unwanted access and data theft. An Intrusion Detection System checks all the network packets and tries to categorize the traffic as disruptive or non-intrusive. The identification of intrusion is the mechanism that starts where the firewall stops.

1.2 WHY DIMENSIONALITY REDUCTION?

Raw Datasets needs to be pre-processed for any type of Intrusion Detection. In this context, reduction of number of dimensions in the high dimensional data is very important. Numerous dimensionality reduction techniques have been used in the last few decades to filter the data samples of the dataset considered. Dimensionality reduction involves the mapping of high-dimensional to lower-dimensionality inputs such that identical points are mapped to neighboring points on the n- manifold in the input space. Dimensionality Reduction decreases a huge data processing overhead on the Machine Learning algorithms and decreases the chances of model being overfitted. This also increases the performance of the algorithm at a very significant rate.

2. OBJECTIVE

Analyze the performance of Machine Learning models with and without Dimensionality Reduction. Reduction is done using two major techniques namely PCA and LDA. Raw Intrusion Detection Data is transformed and decomposed using Label and One Hot Encoders, then this undergoes dimensionality reduction process. The reduced features are then fed to the ML Techniques like DT, NB, SVM and RF.

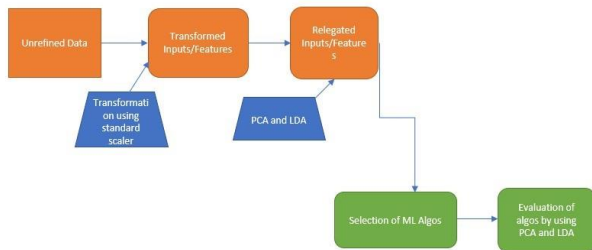


Fig -1: Evaluation of Dimensionality Reduction

In the above figure illustrates the process how the performance evaluation of Dimensionality reduction is done. First, we take the necessary data attributes(features) and labels (customer transactional data) from the dataset and encode categorical data to numerical encoded data. Then perform PCA and LDA to decrease dimensionality of the data.

3. METHODOLOGY

In this project we have used two popular DR algorithms namely Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA) with ML techniques like DT, NB, SVM and RF. Dataset undergoes Data pre-processing, Feature engineering and then application of various machine learning algorithms.

1.1 DATA PRE-PROCESSING

In this process two popular encoding techniques are used namely One-Hot Encoding and Label Encoding. Both these methods are used to deal with categorical data but in a different way. One-Hot encoders simply create a new feature in the Dataset while Label Encoders label them with numerical values.

1.2 FEATURE ENGINEERING

Feature engineering is a very important step in pre-processing of data that aims to remove transformed features from the raw dataset, simplifying the ML model and enhancing the consistency of a machine learning algorithm's output. Machine Learning professionals invest much of their time in data cleaning and engineering functionality. The two techniques used are PCA and LDA.

1.3 PRINCIPAL COMPONENT ANALYSIS

This is a mathematical model that uses orthogonal conversion. A group of correlated variables is transformed by PCA to a group of uncorrelated variables. For exploratory data processing, PCA is used. For analysis of the relationships between a group of variables, PCA may also be used. It can thus be used to minimize dimensionality.

Let us consider a n-dimensional input.

$$x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}$$

The dimensionality is reduced to k -dimensions where $k \ll n$ This is done using following these steps of PCA.

1. Standardization of Data: The data fed must have zero mean and unit variance.

$$x_j^i = \frac{x_j^i - \bar{x}_j}{\sigma_j} \quad \forall_j$$

2. Calculation of the co-variance matrix for the given data.

$$\sum = \frac{1}{m} \sum_i (x_i) (x_i)^T, \quad \sum \in R^{n \times n}$$

3. Calculation of Eigen Vectors and Eigen Values

$$u^T \sum = \gamma u$$

$$U = \begin{pmatrix} | & | & | \\ u_1 & u_2 \dots & u_n \\ | & | & | \end{pmatrix}, u_i \in R^n$$

4. Projection of given raw data as a k-dimensional subspace Now we can choose top k eigen vectors from this matrix formed.

$$x_i^{new} = \begin{pmatrix} u_1^T x^i \\ u_2^T x^i \\ \vdots \\ u_k^T x^i \end{pmatrix} \in R^k$$

1.4 LINEAR DISCRIMINANT ANALYSIS

Another prominent technique for dimensionality reduction is LDA. LDA focuses on reduction of high dimensional data into lower dimensions with a good class level separability which will indeed reduce computation costs. Process of LDA is similar to PCA where PCA uses technique of maximizing variance while LDA maximizes separation of classes. LDA ensures reduction of higher dimensional data into lower subspace say i (where $i \leq x - 1$) without any disturbance in information of class.

BIOGRAPHIES

Venkateswara Gupta Pola pursuing Bachelor of Technology in Computer Science and Engineering at SASTRA University, Thanjavur and working on VQA and Deep learning Research projects at University Autonoma De Barcelona, Spain.



Pranahith Babu Yarra pursuing Bachelor of Technology in Mechatronics Engineering at Mahatma Gandhi Institute of Technology, Hyderabad, currently doing Project on Automatic CNC Laser cutting in MGIT.



Ramya. A pursuing Bachelor of Technology in Computer Science and Engineering at Marri Laxman Reddy Institute of Technology, Hyderabad, currently working on YOLO based blood cells detection in MLRIT.



Sai Suraj Karra pursuing Bachelor of Technology in Mechatronics Engineering at Mahatma Gandhi Institute of Technology, Hyderabad, currently doing Project on Automatic CNC Laser cutting in MGIT.