

Received July 8, 2019, accepted July 23, 2019, date of publication July 30, 2019, date of current version August 14, 2019.

Digital Object Identifier 10.1109/ACCESS.2019.2931956

A Linear Model Based on Principal Component Analysis for Disease Prediction

H. ROOPA¹ AND T. ASHA²

¹Department of Information Science and Engineering, Bangalore Institute of Technology, Bengaluru 560004, India

²Department of Computer Science and Engineering, Bangalore Institute of Technology, Bengaluru 560004, India

Corresponding author: H. Roopa (roopatejas@gmail.com)

ABSTRACT Various classification methods are applied to predict different diseases, such as diabetes, tuberculosis, and so on, in medical field. Diagnosis of diabetes can be analyzed by checking the level of blood sugar of patient with the normal known levels, blood pressure, BMI, skin thickness, and so on. Several classification methods have been implemented on diabetes. In this paper, the main aim is to build a statistical model for diabetes data to get better classification accuracy. Extracted features of diabetes data are projected to a new space using principal component analysis, then, it is modeled by applying linear regression method on these newly formed attributes. The accuracy obtained by this method is 82.1% for predicting diabetes which has reformed over other existing classification methods.

INDEX TERMS Principal component analysis, linear regression model, diabetes.

I. INTRODUCTION

Analysis of diseases is a tough task in medical field. Diagnosing diabetes [2] can be understood by checking the blood sugar level with normal desired level. In this manuscript we present a statistical diabetes disease prediction model where Principal Component Analysis (PCA) is applied to extract attributes of Pima Indian Diabetes Data (PIDD) to a new feature space. These attributes are then modeled using linear regression model [8] to predict diabetes.

Attributes of PIDD are inspected at different angles to obtain required information for processing data. So, feature extraction is a major step in examining PIDD. The work concentrates on retrieving feature values from PIDD to a new feature space by employing PCA method. These new set of feature values are inspected for their importance and relevance, and are subjected for data mining methods like LRM to classify the given data for predicting diabetes disease.

PCA is reduction method which considers the PIDD as set of rows representing characteristics in a high dimensional space and all rows are put up to a directions which represents the best set of features. Here for original set of attributes of PIDD, this transformation is applied to obtain an axis that contains principle eigenvector where all the points of all observations of each feature are spread out. Maximized variance of data can be found on this axis. When considering

The associate editor coordinating the review of this manuscript and approving it for publication was Yue Zhang.

second eigenvector, observe that axis along which the variance of distance from first axis is greatest and so on. Small number of eigenvectors is represented by matrix of points, then to minimize the root mean square error, approximation of data is performed for the given number of columns in the matrix consider. Thus the original features of PIDD are approximated with fewer dimensions which are an overall of original PIDD.

Model building provides a good fit to any set of data. Linear statistical model estimate the unknown dependent PIDD feature value from the known independent PIDD feature values. The representation of relationship between dependent PIDD image feature and set of independent multiple PIDD features are known as regression analysis.

The work is explained as follows, section 2 describes about previous work on PIDD, proposed methodology is explained in section 3, section 4 discusses about findings and section 5 gives about work's conclusion.

II. RELATED WORK

Polat *et al.* [1] proposed a Least Square Support Vector Machine (LS-SVM) classification method to obtain an accuracy of 79.16% which improvised over previous classification methods. Generalized Discriminant Analysis (GDA) was used at preprocessing stage for discriminating variables of PIDD and then LS-SVM technique was applied on these variables for classifying the disease.

