# Review of: "Feature Selection and Classification of Type II Diabetes on High Dimensional Dataset"

Miloš Milovančević

Review of "Feature Selection and Classification of Type II Diabetes on High Dimensional Dataset"

## Abstract

The paper by Priya Vinoth presents a study on feature selection and classification of Type II Diabetes using high-dimensional datasets. The focus is on utilizing the Naïve Bayesian classifier and assessing its performance with different subsets of features from the Pima Indian Type II Diabetes dataset. The study highlights the importance of feature selection in improving classifier performance and reducing computational complexity.

## Introduction

The introduction effectively sets the stage by explaining the significance of feature selection in machine learning and data mining. It describes how feature selection helps in reducing overfitting, improving accuracy, and decreasing training time. The importance of this study in the context of high-dimensional datasets is well articulated, establishing a clear rationale for the research.

## Naive Bayes Classifier

The section on the Naïve Bayes classifier provides a concise explanation of the algorithm. It covers the basic principles, including the use of Bayes' theorem and the assumption of feature independence. The explanation is clear and easy to follow, making it accessible to readers with a basic understanding of probabilistic classifiers.

## Proposed Architecture

The proposed architecture for the study is outlined with a focus on evaluating the Naïve Bayes classifier using different subsets of features. The use of Python and its libraries for machine learning, such as scipy and sklearn, is justified due to their flexibility and extensive functionality. The methodology for feature selection using the SelectKBest class is explained, which provides a practical approach to reducing dataset dimensionality.

## Classification Algorithms Used in this Work

Besides the Naïve Bayes classifier, the study mentions the use of other classification algorithms, including Support Vector

Machine (SVM), Random Forest (RF), Logistic Regression (LR), K-Nearest Neighbor (KNN), and Gradient Boosting Classifier (GBC). However, the paper primarily focuses on the performance of the Naïve Bayes classifier with different feature subsets.

## Methodology Followed

The methodology section is detailed and provides a step-by-step explanation of the processes involved, from loading libraries and data to performing feature selection and standardization. The use of the Pima Indian Diabetes dataset and the steps taken to preprocess and analyze the data are clearly described. The paper employs K-fold cross-validation to evaluate the model, which is a robust method for performance assessment.

## Results

The results demonstrate the impact of feature selection on the performance of the Naïve Bayes classifier. The study finds that a subset of four features yields the best performance in terms of accuracy, precision, recall, and F-measure. This finding is significant as it highlights the effectiveness of feature selection in improving classifier performance while reducing computational requirements.

## Conclusion

The conclusion summarizes the key findings and emphasizes the importance of feature selection in handling high-dimensional datasets. The study's observation that the Naïve Bayes classifier performs optimally with a specific subset of features is an important contribution to the field. The discussion on the "Curse of Dimensionality" and its impact on classifier performance provides valuable insights.

## Overall Evaluation

Priya Vinoth's paper offers a comprehensive analysis of feature selection and its impact on the classification of Type II Diabetes using high-dimensional datasets. The methodology is sound, and the findings are well-supported by the results. The paper is well-organized and provides a clear explanation of the concepts and techniques used.

However, the paper could benefit from a more in-depth comparison of the Naïve Bayes classifier with the other mentioned classifiers. Additionally, a discussion on the potential limitations of the study and suggestions for future research could enhance the paper's contribution to the field. Overall, this study is a valuable addition to the literature on feature selection and machine learning in the context of healthcare data.