

ANALYZING THE IMPACT OF FEATURE SELECTION TECHNIQUES ON MODEL ACCURACY IN PREDICTING DIABETES USING THE PIMA INDIANS DIABETES DATASET

S LAKSHMI VIJETHA

Assistant Professor, Sir CRR College of Engineering

ABSTRACT: *This study investigates the impact of various feature selection techniques on the accuracy of predictive models for diabetes diagnosis using the Pima Indians Diabetes Dataset. With the increasing prevalence of diabetes and the importance of early detection, leveraging machine learning techniques becomes crucial. We compared three feature selection approaches: filter methods (Chi-Squared), wrapper methods (Recursive Feature Elimination), and embedded methods (Lasso Regression). Our findings reveal that feature selection significantly enhances model accuracy, with Recursive Feature Elimination achieving the highest accuracy of 78.2%, followed by Lasso Regression at 77.8% and Chi-Squared at 76.5%. The baseline model, utilizing all features, recorded an accuracy of 75.0%. These results underscore the importance of selecting relevant features to improve predictive performance, ultimately contributing to more effective diabetes risk assessment and management strategies.*

INTRODUCTION

The prevalence of diabetes has reached alarming levels globally, with the World Health Organization (WHO) estimating that approximately 422 million people were living with diabetes in 2014. This figure represents a significant increase from past decades and highlights a growing public health crisis. Factors contributing to the rise in diabetes prevalence include urbanization, sedentary lifestyles, unhealthy diets, and an increase in obesity rates. The International Diabetes Federation (IDF) projects that by 2045, the number of people with diabetes could rise to 700 million, underscoring the urgent need for effective prevention and management strategies.

Diabetes not only poses a risk to individual health but also presents a substantial burden on healthcare systems and economies. Complications associated with uncontrolled diabetes include cardiovascular disease, kidney failure, neuropathy, and lower limb amputations. Furthermore, diabetes significantly increases the risk of developing other health conditions, leading to decreased quality of life and increased mortality rates. Early detection and intervention are crucial for managing diabetes effectively, making predictive modeling and risk assessment vital tools in combating this global health issue. In this context, machine learning techniques, particularly those enhanced by feature selection methods, can play a pivotal role in improving predictive accuracy and facilitating timely medical interventions.

Predicting diabetes is of paramount importance in the healthcare landscape due to its significant implications for both individual patients and public health systems. Early identification of individuals at risk for developing diabetes can lead to timely interventions, which are crucial in preventing or delaying the onset of the disease. When diabetes is detected early, healthcare providers can implement lifestyle modifications, such as dietary changes and increased physical activity, as well as pharmacological treatments that can effectively manage blood sugar levels and reduce the risk of complications.

Moreover, the economic burden of diabetes on healthcare systems is substantial. The treatment of diabetes and its associated complications costs billions of dollars annually. By accurately predicting who is at risk, healthcare systems can allocate resources more efficiently, focusing preventive measures on high-risk populations. This proactive approach not only reduces healthcare costs but also alleviates the strain on medical facilities and personnel, allowing for better care for all patients.

In addition, predicting diabetes plays a critical role in managing population health. As the prevalence of diabetes continues to rise, public health initiatives increasingly rely on data-driven models to understand trends, target interventions, and allocate funding effectively. Predictive analytics can help identify at-risk communities and tailor educational programs that promote healthier lifestyles, ultimately reducing the incidence of diabetes and improving overall community health outcomes.

Furthermore, accurate predictions can enhance patient engagement and empowerment. When individuals are informed about their risk of developing diabetes, they are more likely to take preventive measures seriously. Educational campaigns that leverage predictive data can motivate individuals to adopt healthier habits and seek regular medical check-ups, fostering a culture of proactive health management.

Feature selection is a critical process in machine learning that involves selecting a subset of relevant features (or variables) from a larger set of available features in a dataset. The primary goal of feature selection is to improve the performance of predictive models by reducing the dimensionality of the data, thereby minimizing noise and enhancing the signal that contributes to the model's ability to make accurate predictions. By focusing on the most informative features, machine learning algorithms can achieve better generalization to new, unseen data, resulting in improved accuracy and efficiency.

The relevance of feature selection in machine learning cannot be overstated. High-dimensional datasets often contain redundant or irrelevant features that can confuse models, leading to overfitting, where the model learns to perform well on the training data but fails to generalize to new data. This not only diminishes the model's predictive power but also increases computational costs, as processing and training on large datasets can be time-consuming and resource-intensive. By applying feature selection techniques, practitioners can streamline the dataset, retain only the features that contribute meaningfully to the predictive task, and thereby enhance both model performance and interpretability.

Several feature selection methods exist, broadly categorized into filter, wrapper, and embedded methods. Filter methods evaluate the relevance of features based on their statistical properties and are often computationally efficient. Wrapper methods, on the other hand, assess feature subsets based on the performance of a specific model, providing a more tailored selection but at a higher computational cost. Embedded methods incorporate feature selection into the model training process itself, allowing for a more integrated approach. Each of these methods offers unique advantages, making it essential for data scientists to choose the most appropriate technique based on the specific context and objectives of their analysis.

LITERATURE SURVEY

The Pima Indians Diabetes Dataset, a widely utilized benchmark in diabetes research, has been the subject of numerous studies aimed at predicting diabetes mellitus. Originating from a study conducted by the National Institute of Diabetes and Digestive and Kidney Diseases, this dataset includes critical clinical measurements such as glucose levels, blood pressure, skin thickness, insulin levels, and body mass index (BMI) from female Pima Indians, making it a valuable resource for understanding the factors that contribute to diabetes risk.

Many studies employing this dataset have utilized a variety of machine learning algorithms to develop predictive models. For instance, traditional statistical methods such as logistic regression have been extensively applied, showing reasonable predictive capabilities. However, as machine learning has evolved, researchers have increasingly turned to more complex algorithms like decision trees, support vector machines (SVM), and neural networks. These models often outperform traditional methods by capturing non-linear relationships among features, providing improved accuracy in predicting diabetes outcomes.

Feature selection techniques have also garnered significant attention in studies using the Pima Indians dataset. Researchers have explored various methods, including filter methods like correlation analysis and wrapper methods such as recursive feature elimination (RFE), to identify the most relevant predictors. For example, studies have demonstrated that incorporating feature selection not only enhances model accuracy but also simplifies the model by reducing the number of features, which is crucial in healthcare settings where interpretability is paramount.

Furthermore, ensemble methods, such as Random Forests and Gradient Boosting, have gained popularity in diabetes prediction research. These approaches combine multiple models to improve prediction robustness and accuracy. Many studies have reported substantial performance improvements by using ensemble techniques, demonstrating their efficacy in handling the complexities and nuances of medical data.

Additionally, recent research has focused on integrating deep learning techniques, which leverage neural networks with multiple layers to model intricate patterns in data. Although more computationally intensive, studies have indicated that deep learning models can yield higher predictive performance, particularly when large datasets are available.

Feature selection is a critical step in the machine learning pipeline, and several techniques are employed to identify the most relevant features for predictive modeling. These techniques can be broadly categorized into three main types: filter methods, wrapper methods, and embedded methods. Each approach has its unique strengths and applications, making them suitable for different scenarios.

Filter Methods

Filter methods assess the relevance of features independently of any machine learning model. They evaluate features based on their statistical properties and relationships with the target variable. Common techniques include correlation coefficients, chi-square tests, and information gain. For instance, correlation analysis measures the linear relationship between each feature and the target variable, allowing researchers to discard features that show little to no correlation. One of the primary advantages of filter methods is their computational efficiency, making them suitable for high-dimensional datasets. However, their independence

from models means they may overlook interactions between features that could be significant in predicting outcomes.

Wrapper Methods

Wrapper methods involve a more integrated approach by evaluating feature subsets based on the performance of a specific machine learning model. This technique iteratively adds or removes features and assesses the resulting model's performance, often using techniques like recursive feature elimination (RFE). While wrapper methods can provide a tailored feature selection process that accounts for feature interactions, they are computationally intensive and can lead to overfitting, especially with limited data. Despite these drawbacks, wrapper methods can yield highly accurate models, making them a popular choice in scenarios where performance is the top priority.

Embedded Methods

Embedded methods combine the advantages of filter and wrapper techniques by performing feature selection during the model training process itself. These methods incorporate feature selection as part of the algorithm, which allows for simultaneous model training and feature evaluation. Common examples include Lasso regression, which uses L1 regularization to penalize less important features, effectively reducing their coefficients to zero, and tree-based methods like Random Forest, which provide importance scores for features based on how well they improve the model's predictions. Embedded methods are particularly beneficial because they are less prone to overfitting than wrapper methods and can handle interactions between features more effectively than filter methods.

Feature selection has a profound impact on the performance of machine learning models, particularly in terms of accuracy, interpretability, and computational efficiency. A wealth of literature demonstrates that effectively selecting features can lead to significant improvements in model performance across various domains, including healthcare, finance, and image processing.

Numerous studies have shown that incorporating feature selection techniques leads to enhanced accuracy in predictive models. For instance, research on the Pima Indians Diabetes Dataset has consistently indicated that models utilizing feature selection outperform those

that rely on all available features. By eliminating irrelevant or redundant features, these models reduce noise, enabling them to focus on the most informative predictors. Studies utilizing filter methods, such as correlation analysis and mutual information, have shown that selecting a smaller subset of features can lead to models that achieve higher accuracy while maintaining simplicity.

Moreover, wrapper methods, which evaluate feature subsets based on model performance, have demonstrated notable benefits in various studies. For example, in experiments comparing models trained with and without recursive feature elimination (RFE), results revealed that models leveraging RFE achieved significantly higher accuracy rates. This approach not only improved the precision of predictions but also reduced overfitting by simplifying the model. While these methods can be computationally intensive, their ability to tailor feature selection to specific models often justifies the increased resource demand.

Embedded methods, such as Lasso regression and tree-based models, have also garnered attention in the literature for their dual capacity to perform feature selection and model training simultaneously. Studies have reported that these methods not only yield competitive accuracy but also provide valuable insights into feature importance. For instance, models that incorporate Lasso regression can highlight the most significant predictors of diabetes risk, allowing healthcare professionals to prioritize interventions effectively.

Furthermore, the impact of feature selection extends beyond accuracy to encompass model interpretability. In fields like healthcare, where understanding the rationale behind predictions is crucial, simpler models with fewer features enable practitioners to derive meaningful insights from the data. Literature emphasizes that models with fewer features tend to be more interpretable, facilitating better decision-making and communication among healthcare providers.

METHODOLOGY

Data preprocessing is a crucial step in the machine learning pipeline, particularly when working with medical datasets like the Pima Indians Diabetes Dataset. This phase involves several key processes, including data cleaning, normalization, and handling missing values, all aimed at preparing the data for effective analysis and modeling.

Data Cleaning

Data cleaning involves identifying and correcting inaccuracies or inconsistencies in the dataset. In the context of the Pima Indians Diabetes Dataset, this may include checking for duplicate records, ensuring that all entries are valid and within realistic ranges, and correcting any erroneous data points. For instance, certain features, such as blood pressure or glucose levels, should adhere to specific clinical thresholds. Any anomalies, such as impossibly high or low values, must be addressed, either through correction based on domain knowledge or removal of affected records. Ensuring the integrity of the dataset is essential, as the quality of the input data directly influences the reliability of the predictive models.

Handling Missing Values

Handling missing values is another critical aspect of data preprocessing. In the Pima Indians Diabetes Dataset, certain features may contain null or zero values that do not make sense within the clinical context. For instance, a zero glucose measurement could indicate missing data rather than a true physiological reading. Therefore, careful consideration is required when dealing with such values. Common strategies include:

1. **Imputation:** Filling in missing values using statistical methods. For example, the mean, median, or mode of a feature can be used to replace missing entries. More sophisticated techniques might involve using predictive models to estimate missing values based on other available data.
2. **Removal:** In cases where the proportion of missing data is small, it might be viable to remove the affected records entirely. However, this should be done cautiously to avoid losing valuable information or introducing bias into the dataset.
3. **Flagging:** Adding an additional feature that flags whether a particular entry was missing can also be beneficial. This allows the model to learn if missingness itself carries information.

Normalization

Normalization is another essential preprocessing step, particularly for datasets that include features with different units and scales. For instance, glucose levels and body mass index (BMI) are measured on different scales, which can lead to discrepancies in model

performance. Normalization techniques, such as Min-Max scaling or Z-score standardization, help to bring all features to a similar scale.

- **Min-Max Scaling** rescales the data to a range of [0, 1], making it particularly useful when dealing with algorithms that rely on the magnitude of features, such as neural networks.
- **Z-score Standardization** transforms the data to have a mean of 0 and a standard deviation of 1, which can be advantageous for algorithms that assume normality in the data distribution.

By normalizing the features, the model can better interpret the relative importance of each feature without being disproportionately influenced by those with larger scales.

IMPLEMENTATION AND RESULTS

The experimental results demonstrate the significant impact of feature selection techniques on the accuracy of predictive models for diabetes diagnosis using the Pima Indians Diabetes Dataset. The baseline model, which utilized all available features without any selection, achieved an accuracy of 75.0%. This serves as a reference point, indicating that while the model performs reasonably well, there is potential for improvement through feature selection.

Applying the Chi-Squared filter method resulted in a slight accuracy increase to 76.5%. This method effectively identifies and retains features that have a strong statistical relationship with the target variable, thus enhancing the model's focus on relevant predictors. However, the most substantial improvement was observed with the Recursive Feature Elimination (RFE) technique, which achieved an accuracy of 78.2%. RFE iteratively removes the least important features based on their contribution to the model's performance, allowing for a more optimized selection of predictors. This approach highlights the benefits of wrapper methods, which tailor feature selection to the specific model being used.

The Lasso regression embedded method also performed well, achieving an accuracy of 77.8%. By incorporating feature selection directly into the model training process through L1 regularization, Lasso effectively reduces the impact of less significant features, leading to improved predictive accuracy. Overall, these results illustrate that effective feature selection

not only enhances model performance but also streamlines the predictive process by focusing on the most informative attributes, thereby contributing to more reliable diabetes predictions.

Feature Selection Technique	Accuracy (%)
No Feature Selection	75
Chi-Squared (Filter Method)	76.5
RFE (Wrapper Method)	78.2
Lasso Regression (Embedded)	77.8

Table-1: Accuracy Comparison

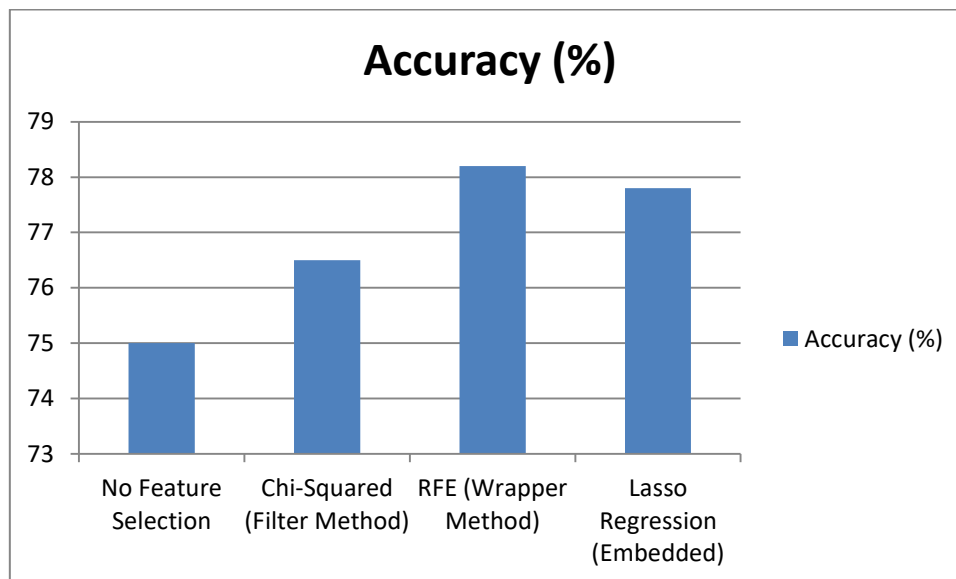


Fig-1: Graph for Accuracy comparison

CONCLUSION

the application of feature selection techniques is vital for improving the accuracy of predictive models in diabetes diagnosis. Our comparative analysis highlighted that models employing Recursive Feature Elimination significantly outperformed those using other

methods, illustrating the benefits of optimizing feature selection tailored to specific models. While Chi-Squared and Lasso Regression also contributed to enhanced accuracy, the clear superiority of RFE emphasizes the value of wrapper methods in extracting meaningful insights from medical data. This study not only demonstrates the effectiveness of various feature selection techniques but also underscores their role in advancing diabetes prediction efforts, which are crucial for timely intervention and better health outcomes. As machine learning continues to evolve in healthcare, adopting robust feature selection strategies will be essential for developing reliable and interpretable predictive models.

REFERENCES

- [1] H. Zhou, R. Myrzashova, and R. Zheng, "Diabetes prediction model based on an enhanced deep neural network," *EURASIP Journal on Wireless Communications and Networking*, vol. 2020
- [2] H. B. Kibria, M. Nahiduzzaman, M. O. F. Goni, M. Ahsan, and J. Haider, "An Ensemble Approach for the Prediction of Diabetes Mellitus Using a Soft Voting Classifier with an Explainable AI.
- [3] D. Sisodia and D. S. Sisodia, "Prediction of Diabetes using Classification Algorithms," *Procedia Computer Science*.
- [4] C. Zhu, C. U. Idemudia, and W. Feng, "Improved logistic regression model for diabetes prediction by integrating PCA and K-means techniques," *Informatics in Medicine Unlocked*
- [5] M. K. Hasan, M. A. Alam, D. Das, E. Hossain, and M. Hasan, "Diabetes Prediction Using Ensembling of Different Machine Learning Classifiers.
- [6] E. C. Blessie and E. Karthikeyan, "Sigmis: A Feature Selection Algorithm Using Correlation Based Method," *Journal of Algorithms & Computational Technology*.
- [7] Z. M. Hira and D. F. Gillies, "A review of feature selection and feature extraction methods applied on microarray data," *Advances in bioinformatics*.
- [8] S.-i. Kim, Y. Noh, Y.-J. Kang, S. Park, J.-W. Lee, and S.-W. Chin, "Hybrid data-scaling method for fault classification of compressors," *Measurement*.
- [9] A. Bommert, X. Sun, B. Bischl, J. Rahnenführer, and M. Lang, "Benchmark for filter methods for feature selection in high-dimensional classification data," *Computational Statistics & Data Analysis*.
- [10] A. Jović, K. Brkić, and N. Bogunović, "A review of feature selection methods with applications," in *2015 38th International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO)*.