

Data Analytics Lab: Assignment-6

A Mathematical Essay on Support Vector Machine

Arjav Singh
Metallurgical and Materials Engineering
Indian Institute of Technology Madras
Chennai, India
mm20b007@smail.iitm.ac.in

Abstract—This study is a mathematical exposition of the Support Vector Machine method, which is applied to a dataset containing different star attributes. The primary objective is to differentiate pulsar stars from non-pulsar stars by utilizing different metrics derived from their integrated pulse profiles (folded profiles). The essay aims to assess the efficacy of SVMs in this classification task and compare the performance of different kernels. Additionally, it investigates whether specific features exhibit distinctive characteristics crucial for accurate classification between the two-star classes.

Index Terms—Introduction, Support Vector Machine, Data & Problem, Conclusion

I. INTRODUCTION

Pulsars are rotating neutron stars observed to have pulses of radiation at regular intervals that typically range from milliseconds to seconds. Pulsars have very strong magnetic fields which funnel jets of particles out along the two magnetic poles. These accelerated particles produce very powerful beams of light.

This study focuses on a comprehensive empirical classification of a pulsar from a normal star based on its features, including the *Mean of the integrated profile*, *Excess kurtosis of the integrated profile*, *Skewness of the integrated profile*, *Mean of the DM-SNR curve*, *Excess kurtosis of the DM-SNR curve*, *Skewness of the DM-SNR curve*. Support Vector Machine Classification Machine learning technique is used to achieve the goal. In SVM, the data points are first represented in an n-dimensional space. The algorithm then uses statistical approaches to find the best line that separates the various classes present in the data.

The research methodology initiates with the acquisition, refinement, and preprocessing of the raw data. An exploratory data analysis follows this to gain a deeper understanding of the dataset's inherent features. Subsequently, statistical models are crafted, along with the generation of visual aids to offer both quantitative and visual support for the observed associations. The subsequent section furnishes an exposition and discourse on the insights and revelations extracted from the data analysis and the models that have been generated. It emphasizes the significant discoveries, recurring patterns, and emerging trends that have surfaced during the course of the study.

The concluding section summarizes the key highlights and significant features of the research. Potential avenues for further investigation are outlined, suggesting areas where future

research could expand upon the findings. A contribution is made to a deeper understanding of pulsars and how to identify them, with valuable insights for their detection.

II. SUPPORT VECTOR MACHINE

The Support Vector Machine (SVM) is a robust supervised algorithm ideally suited for handling complex datasets. SVM can be employed for regression and classification tasks, although it typically excels in classification problems. Despite being developed in the 1990s, SVM remains a popular choice, known for its high-performance capabilities even with minimal parameter tuning.

A. Types and Features of SVM

- 1) **Linear SVM**: When the data is perfectly linearly separable, we can only use Linear SVM. Perfectly linearly separable means that the data points can be classified into 2 classes by using a single straight line (if 2D).
- 2) **Non-Linear SVM**: When the data is not linearly separable, then we can use Non-Linear SVM, which means when the data points cannot be separated into 2 classes by using a straight line (if 2D) then we use some advanced techniques like kernel tricks to classify them. In most real-world applications, we do not find linearly separable data points; hence we use kernel tricks to solve them.

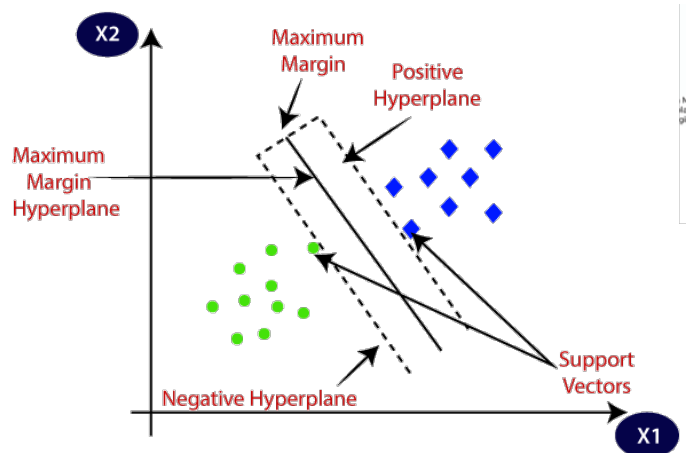


Fig. 1. General Description of Support Vector Machine

An SVM has two major components(Figure 1), which are

- 1) **Support Vectors**:: These are the points closest to the hyperplane. A separating line will be defined with the help of these data points.
- 2) **Margin**: It is the distance between the hyperplane and the observations closest to the hyperplane (support vectors). In SVM large margin is considered a good margin. There are two types of margins hard margin and soft margin.

B. Working of SVM

In SVMs, we mainly aim to select a hyperplane with the maximum possible margin between support vectors in the given dataset. SVM searches for the maximum margin hyperplane in the following 2-step process –

- 1) Generate hyperplanes that segregate the classes in the best possible way. Many hyperplanes might classify the data. We should look for the best hyperplane representing the largest separation, or margin, between the two classes.
- 2) So, we choose the hyperplane so that the distance from it to the support vectors on each side is maximized. If such a hyperplane exists, it is known as the maximum margin hyperplane, and the linear classifier it defines is known as a maximum margin classifier.

Figure 2 illustrates the concept of maximum margin and maximum margin hyperplane in a clear manner.

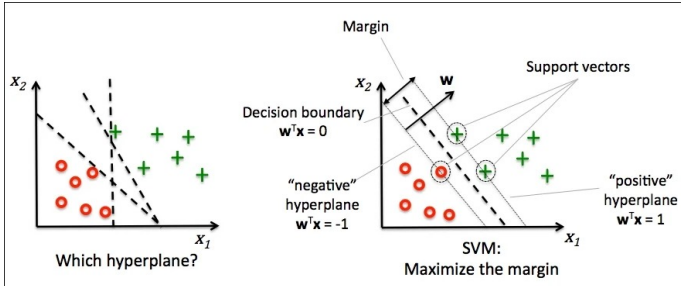


Fig. 2. Working of SVM

C. Problem with dispersed datasets

Sometimes, the sample data points are so dispersed that it is impossible to separate them using a linear hyperplane. In such a situation, SVMs use a kernel trick to transform the input space to a higher dimensional space, as shown in Figure 3. It uses a mapping function to transform the 2-D input space into the 3-D input space. Now, we can easily segregate the data points using linear separation.

D. Kernel Methods for SVM

In practice, SVM algorithm is implemented using a kernel. It uses a technique called the kernel trick. Simply put, a kernel is just a function that maps the data to a higher dimension where data is separable. A kernel transforms a low-dimensional input data space into a higher-dimensional

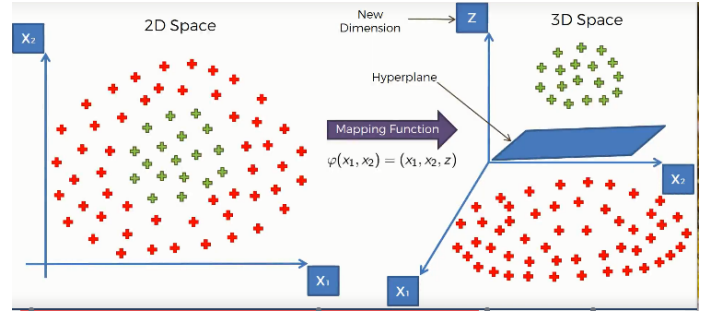


Fig. 3. Working of SVM on non-linear data

space. So, it converts non-linear separable problems to linear separable problems by adding more dimensions to it. Thus, the kernel trick helps us to build a more accurate classifier. Hence, it is useful in non-linear separation problems. We can define a kernel function as follows-

$$K(\bar{x}) = \begin{cases} 1 & \text{if } \|\bar{x}\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Fig. 4. Kernel function

In the context of SVMs, there are 4 popular kernels – Linear kernel, Polynomial kernel, Radial Basis Function (RBF) kernel (also called Gaussian kernel), and Sigmoid kernel. These are described below -

- 1) **Linear kernel**: In linear kernel, the kernel function takes the form of a linear function as follows linear kernel: $K(x_i, x_j) = x_i^T x_j$. Linear kernel is used when the data is linearly separable. It means that data can be separated using a single line. It is one of the most common kernels to be used. It is mostly used when there are large number of features in a dataset. Linear kernel is often used for text classification purposes. Training with a linear kernel is usually faster because we only need to optimize the C regularization parameter. When training with other kernels, we also need to optimize the parameter. So, performing a grid search will usually take more time. The linear kernel can be visualized in Figure 5.
- 2) **Polynomial kernel**: Polynomial kernel represents the similarity of vectors (training samples) in a feature space over polynomials of the original variables. The polynomial kernel looks not only at the given features of input samples to determine their similarity but also at combinations of the input samples. For d-degree polynomials, the polynomial kernel is defined as follows: Polynomial kernel:

$$K(x_i, x_j) = (\gamma x_i^T x_j + r)^d, \gamma > 0$$

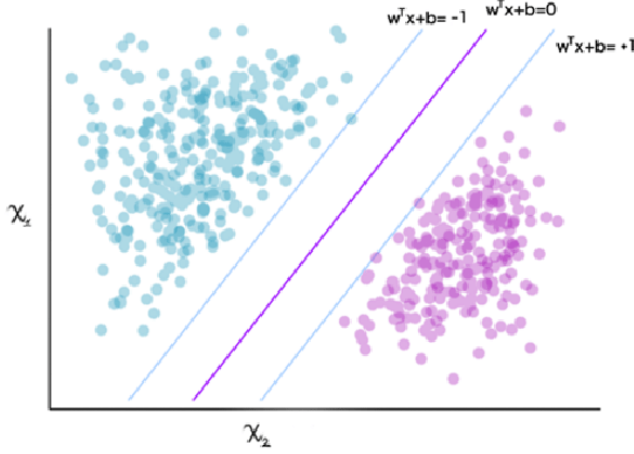


Fig. 5. Linear Kernel function

The polynomial kernel is very popular in Natural Language Processing. The most common degree is $d = 2$ (quadratic) since larger degrees tend to overfit NLP problems. It can be visualized in Figure 6.

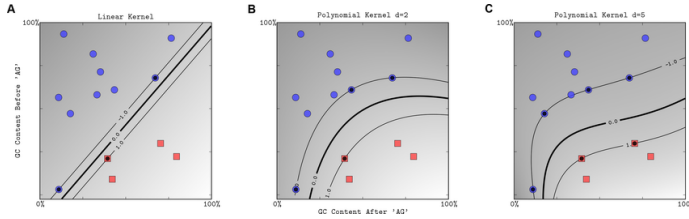


Fig. 6. Polynomial Kernel function

- 3) **Radial basis function kernel:** Radial basis function kernel is a general purpose kernel. It is used when we have no prior knowledge about the data. The RBF kernel on two samples, x and y , is defined by the following equation –

$$k(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$$

Fig. 7. Radial Basis function Kernel Equation

- 4) **Sigmoid Function Kernel:** The sigmoid kernel originates in neural networks and can be used as a proxy for neural networks. The following equation gives the sigmoid kernel: Sigmoid kernel:

$$k(x, y) = \tanh(\alpha x^T y + c)$$

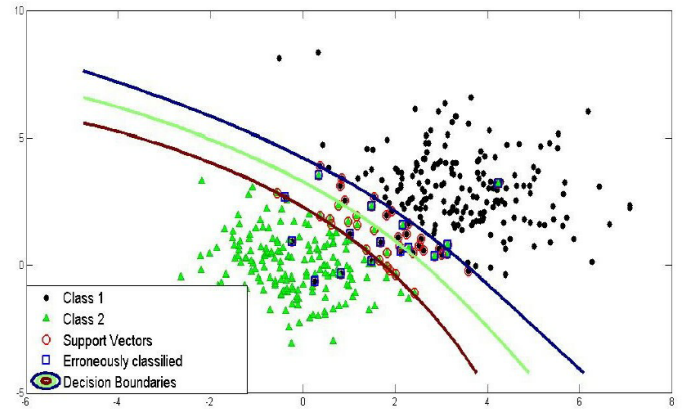


Fig. 8. Classification using radial basis function kernel

E. Metrics for model evaluation

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Fig. 9. Confusion Matrix.

- 1) **Confusion Matrix:** It is used to summarize the performance of a classification algorithm on a set of test data for which the true values are previously known. Sometimes it is also called an error matrix. Terminologies of the Confusion matrix (Figure 1) are:

- **True Positive:** TP means the model predicted yes, and the actual answer is also yes.
- **True negative:** TN means the model predicted no, and the actual answer is also no.
- **False positive:** FP means the model predicted yes, but the actual answer is no.
- **False negative:** FN means the model predicted no, but the actual answer is yes.

The rates calculated using the Confusion Matrix are:

- a) **Accuracy:** $(TP+TN/Total)$ tells about overall how classifier is correct.

- b) **True positive rate:** $TP / (\text{actual yes})$ it says about how much time yes is predicted correctly. It is also called “sensitivity” or “recall.”
- c) **False positive rate:** $FP / (\text{actual number})$ says how much time yes is predicted when the actual answer is no.
- d) **True negative rate:** $TN / (\text{actual number})$ says how much time no is predicted correctly, and the actual answer is also no. It is also known as “specificity.”
- e) **Misclassification rate:** $(FP + FN) / (\text{Total})$ It is also known as the error rate and tells about how often our model is wrong.
- f) **Precision:** $(TP / (\text{predicted yes}))$ If it predicts yes, then how often is it correct.
- g) **Prevalence:** $(\text{actual yes} / \text{total})$ how often yes condition actually occurs.
- h) **F1-score:** f1 score is defined as the weighted harmonic mean of precision and recall. The best achievable F1 score is 1.0, while the worst is 0.0. The F1 score serves as the harmonic mean of precision and recall. Consequently, the F1-score consistently yields lower values than accuracy measures since it incorporates precision and recall in its computation. When evaluating classifier models, it is advisable to employ the weighted average of the F1 score instead of relying solely on global accuracy.

2) **ROC curve (Receiver Operating Characteristic):** The Receiver Operating Characteristic (ROC) curve is a useful tool for assessing a model’s performance by examining the trade-offs between its True Positive (TP) rate, also known as sensitivity, and its False Negative (FN) rate, which is the complement of specificity. This curve visually represents these two parameters. The Area Under the Curve (AUC) metric to summarize the ROC curve concisely. The AUC quantifies the area under the ROC curve. In simpler terms, it measures how well the model can distinguish between positive and negative cases. A higher AUC indicates better classifier performance.

In essence, AUC categorizes model performance as follows:

- If $AUC = 1$, the classifier correctly distinguishes between all the Positive and Negative class points.
- If $0.5 < AUC < 1$, the classifier will distinguish the positive class value from the negative one because it finds more TP and TN than FP and FN.
- If $AUC = 0.5$, the classifier cannot distinguish between positive and negative values.
- If $AUC = 0$, the classifier predicts all positive as negative and negative as positive.

III. PROBLEM

We have been tasked to analyze various attributes of stars, such as their Mean of the integrated profile, Excess kurtosis of the integrated profile, Skewness of the integrated profile, Mean

of the DM-SNR curve, Excess kurtosis of the DM-SNR curve, Skewness of the DM-SNR curve. The goal is to identify which of them are pulsars.

A. Exploratory Data Analysis and Feature Generation

The data is initially read into a pandas data frame. A total of 12528 data points are observed, with 9 columns encompassing various car-related features. When the distributions of the target variable are visualized, a multi-class imbalanced dataset problem is evident. Around 90.8% of the total stars are classified as not pulsars, and only 9.2% are classified as pulsars (as shown in Figure 11). It is observed that 8 out of 9 features are continuous and are numerical. The target feature is categorical in nature and has labels 0 and 1. The aim is to predict this target feature.

Further analysis of the data implies that three features, namely *Excess kurtosis of the integrated profile*, *Standard deviation of the DM-SNR curve*, and *Skewness of the DM-SNR curve* have missing values.

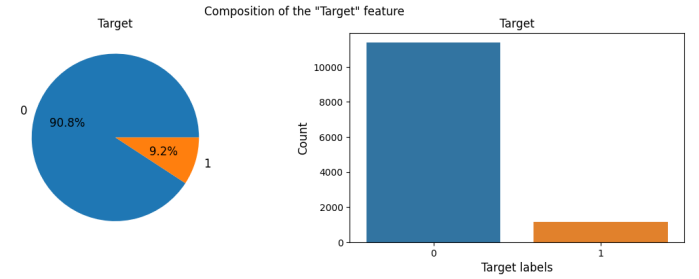


Fig. 10. Distribution of Target Variable

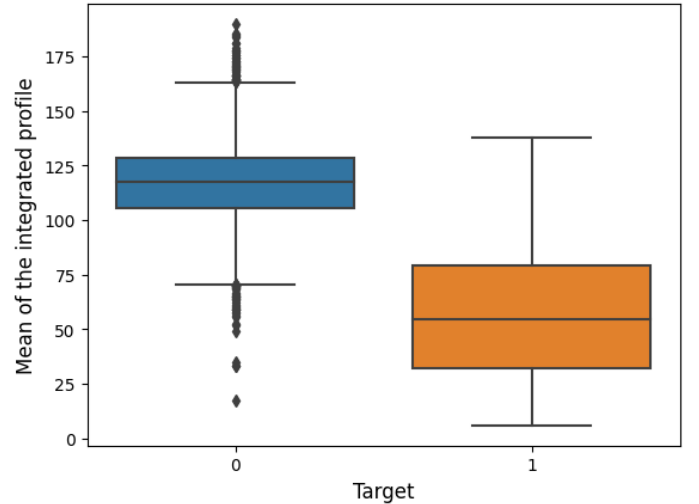


Fig. 11. Mean of the integrated profile vs Target

Univariate analysis is initiated by generating boxplots for each of the eight features, employing the seaborn library, with the target column as the hue. It is observed that there is an opposite trend for the same values of integrated profile and DM-SNR. Correlation is then checked, and Excess kurtosis of

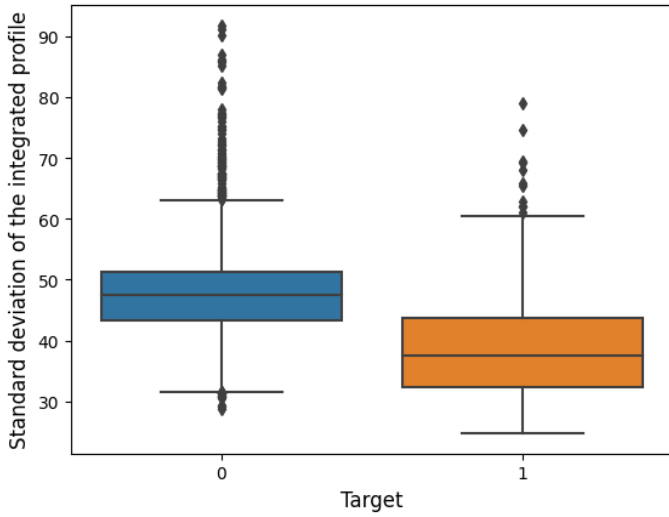


Fig. 12. Standard deviation of the integrated profile vs Target

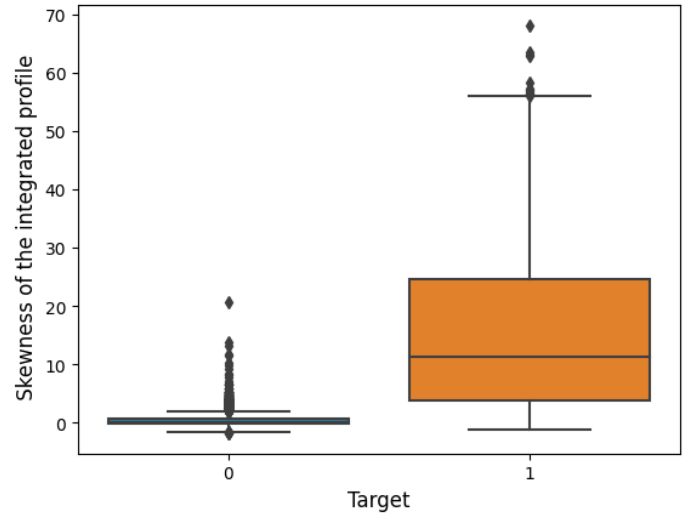


Fig. 14. Skewness of the integrated profile vs Target

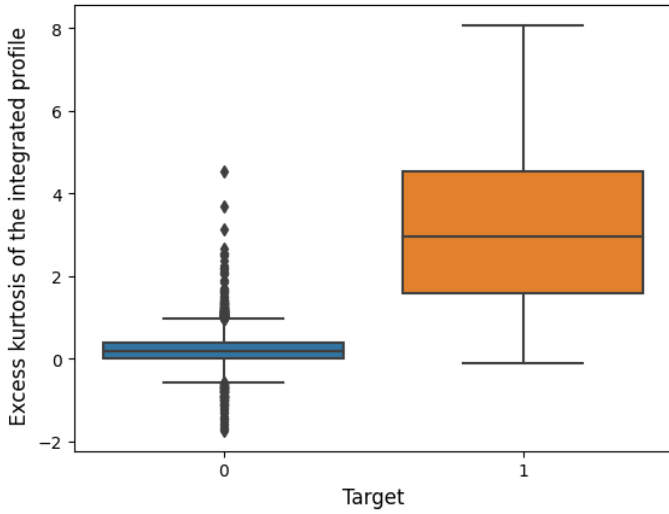


Fig. 13. Excess kurtosis of the integrated profile vs Target

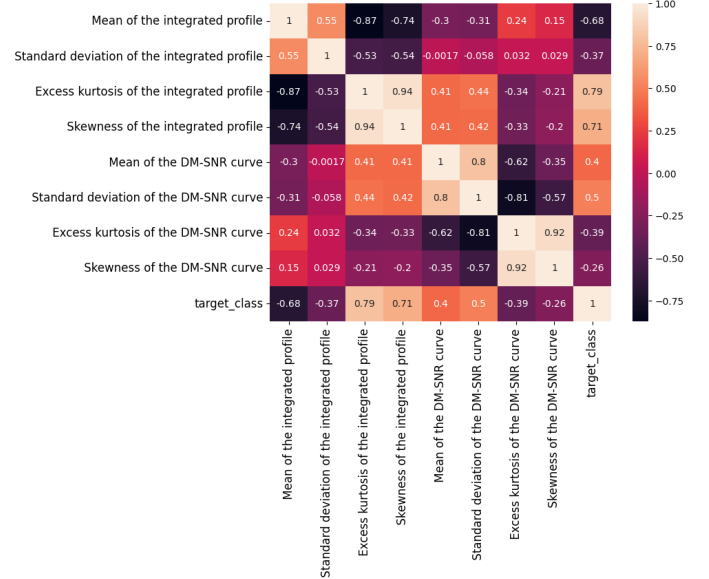


Fig. 15. Correlation heatmap

the integrated profile, Skewness of the integrated profile, and Standard deviation of the DM-SNR curve have a high positive correlation (greater than 0.5) with the target class feature. In contrast, apart from the Mean of the DM-SNR curve, all remaining features negatively associate with the target class feature.

B. Post-Processing and Feature Selection

Since the dataset comprises missing values, it is necessary to handle it properly. There are many methods of imputation, for the given problem, I have used two methods of imputation:

- 1) Standard imputation: The missing values were replaced with median and mean for the respective features.
- 2) Iterative imputation: This method first fits the data and generates a function that predicts the missing values. It

is visible in Figure 17 that this is a better version of the imputation for the given dataset type.

Ultimately, the data is divided using an 80/20 split, resulting in a final dataset with 10022 examples in the training set and 2506 in the cross-validation set. We then use the Standard scaler to scale our train and validation values.

C. SVM Modelling

Modeling is initiated using the default hyperparameters provided by the SVC library in scikit-learn, where the defaults are set as follows: $C = 1.0$, $kernel = rbf$, and $gamma = auto$, among other parameters. Initial observations with default hyperparameters show that an accuracy score of 0.9816, and a

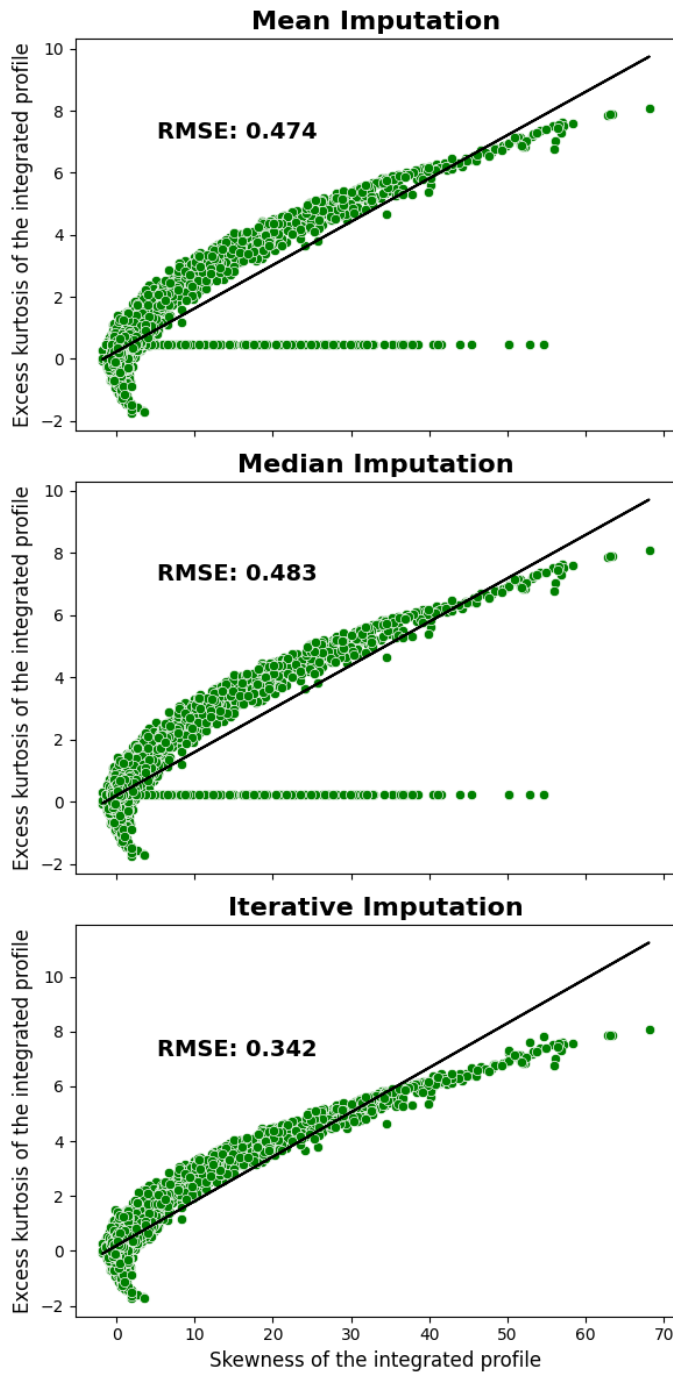


Fig. 16. Simple and Iterative imputation

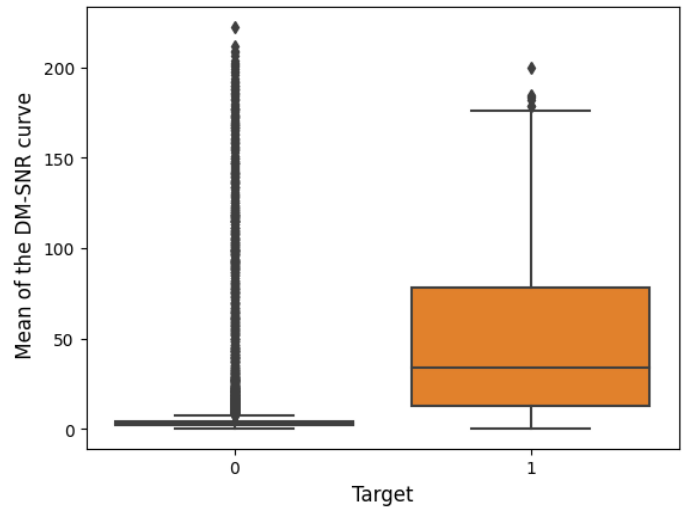


Fig. 17. Mean of the DM-SNR curve vs Target

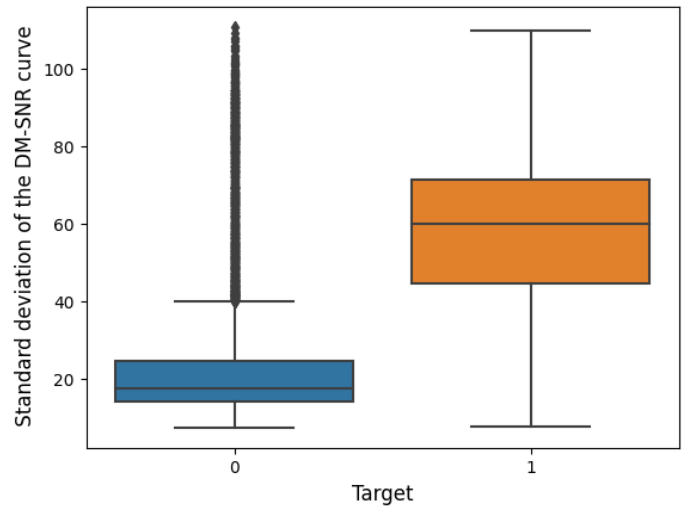


Fig. 18. Standard deviation of the DM-SNR curve vs Target

precision of 0.9594 are achieved by the model. These metrics imply the good performance of the model.

However, it is essential to note that our dataset is imbalanced. In this context, accuracy alone is an inadequate measure for assessing predictive performance. Alternative metrics that offer better insights into model selection must be explored. In particular, attention is turned to the F1 score, which is more informative when dealing with imbalanced datasets. It is found that it is found that the model achieved an F1 score of 0.8915 and an ROC AUC score of 0.9145.

To further enhance model performance, hyperparameter tuning is performed. Grid Search is employed to explore a predefined hyperparameter space, which includes testing various kernels and a range of C values from 1 to 10. Additionally, experimentation is done with the degree for the linear kernel and class weights are set as "balanced." Other kernel methods were also experimented with varying values

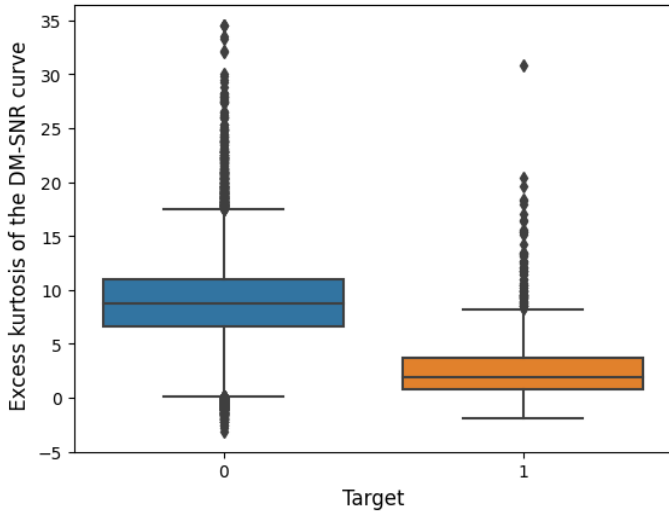


Fig. 19. Excess kurtosis of the DM-SNR curve vs Target

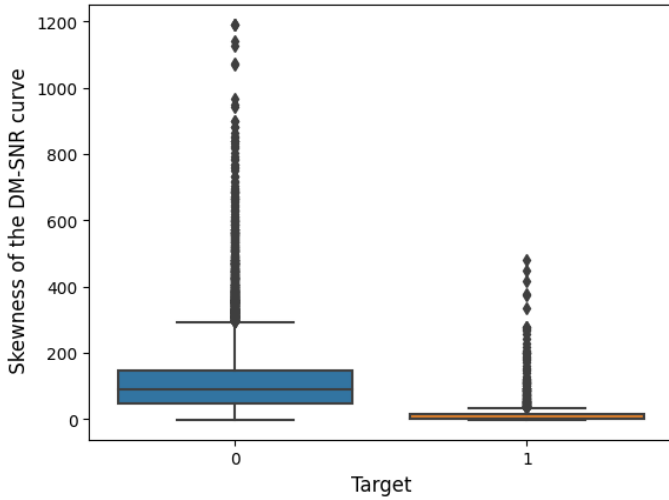


Fig. 20. Skewness of the DM-SNR curve vs Target

of the gamma hyperparameter, for polynomial kernel, different degrees were experimented with, ranging from 2 to 5. Utilizing a 2-fold cross-validation technique.

The best model is identified, with the following parameters ' C ': 10, ' $kernel$ ': *linear* exhibiting an F1 score of 0.9019, an accuracy of 0.9832, a precision of 0.9602, and an ROC AUC score of 0.9234. This represents an improvement over the initial F1 score of 0.8915.

IV. CONCLUSION

Having conducted a comprehensive analysis of the Support Vector methods with various kernels, an improvement of 0.01 in the F1 score was observed using GridSearchCV. Consequently, GridSearchCV serves the purpose of identifying the parameters that will enhance the performance of this specific model. The dataset contains outliers, and as the value of C was increased to reduce the influence of outliers, accuracy

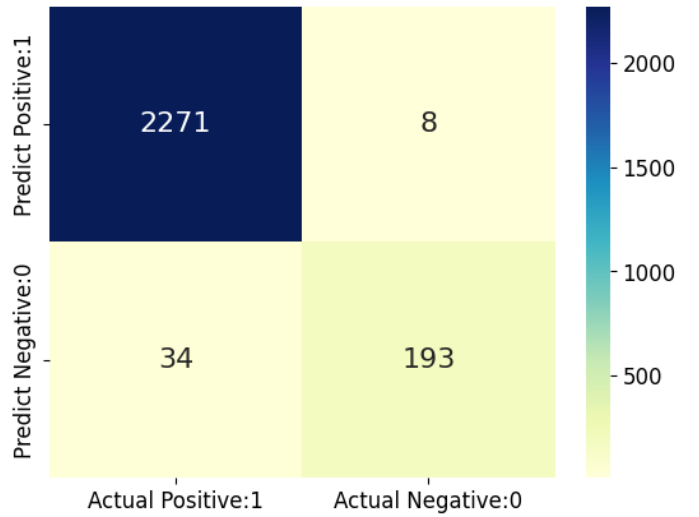


Fig. 21. Confusion Matrix for SVM after grid search

improved. This effect was consistent across different kernel types.

The ROC AUC of the model is very close to 1, suggesting that the classifier excels in classifying pulsar stars. Additionally, the precision and recall values are commendable, standing at 0.9602 and 0.985, respectively. The true positive rate is 0.985, while the false positive rate is 0.191.

In conclusion, SVMs with the linear kernel demonstrate the capability to fit the training data effectively, as anticipated based on the multivariate analysis, which revealed that most features partition the target classes into distinct and easily separable regions. Future possibilities for improvement include exploring additional features that may provide better insights into the target variable. Additionally, addressing feature correlation by eliminating specific features or creating hybrid features could be explored. Given the highly imbalanced data, implementing upsampling and downsampling techniques is another avenue worth considering."

For future work, further avenues of growth could involve exploring additional features that might better explain the target variable.

REFERENCES

- [1] "Everything About SVM Classification: Above and Beyond," Online. Available: <https://towardsdatascience.com/everything-about-svm-classification-above-and-beyond-cc665bfd993e>.
- [2] P. 111, "SVM Classifier Tutorial," Kaggle, Available: <https://www.kaggle.com/code/prashant111/svm-classifier-tutorial/notebook>.
- [3] "Support Vector Machines (SVM): A Complete Guide for Beginners," Analytics Vidhya. [Online]. Available: <https://www.analyticsvidhya.com/blog/2021/10/support-vector-machines-svm-a-complete-guide-for-beginners/>.
- [4] "AUC-ROC Curve & Confusion Matrix Explained in Detail," [Online]. Available: <https://www.kaggle.com/code/vithal2311/auc-roc-curve-confusion-matrix-explained-in-detail>.
- [5] Analytics Vidhya. "K-Fold Cross-Validation Technique and Its Essentials." [Online]. Available: <https://www.analyticsvidhya.com/blog/2022/02/k-fold-cross-validation-technique-and-its-essentials/>.