# Assignment-6: A mathematical essay on Support Vector Machine

Vasudev Gupta

*Interdisciplinary Dual Degree - Data Science*
*Indian Institute of Technology Madras*
Chennai, India
me18b182@smail.iitm.ac.in

*Abstract*—We will present a mathematical research on Support Vector Machines in this paper. We will use the Support Vector Machines to identify whether a star is a pulsar star or not based on the given features and will attempt to determine whether using the Support Vector Machines was the correct decision. We will also attempt exploratory data analysis and feature selection in order to give the correct data to the model. Finally, we'll check if our model can generalise effectively to unseen data.

## I. Introduction

Machine learning techniques are becoming increasingly popular, and they are assisting a number of organisations in tackling difficult challenges using data-driven approaches. Machine learning methods attempt to estimate the function on provided training data by capturing the underlying distribution of the training data and then predict on unseen data using the learned distribution. When building the model, one expects that the training data is comparable to previously unseen data and that the model will generalise well to future data.

SVM maps training examples to points in space to maximise the distance between the two categories. New examples are then mapped into that space and predicted to belong to a class based on where they fall on the margin. SVMs can perform non-linear classification efficiently using the kernel trick, implicitly mapping their inputs into high-dimensional feature spaces, in addition to performing linear classification. An unsupervised learning approach is required, in which the data is naturally clustered into groups, and then new information is mapped to these formed groups.

This paper looks at a geometric model that works on pulsar star recognition data extracted. The key focus of our work is to identify whether a star is a pulsar star or not based on the given features. We visualize and analyze the data to see if any specific elements strongly influence the target variable. Ultimately, we would like to understand the data characteristics which affect the performance of the Support Vector Classifier.

This paper is structured as follows. In the next section, we provide the necessary background and definitions. Section II focuses on the data description, analysis, handling of missing values and visualization to gain more insights into the problem. At the same time, Section III demonstrates that the use of Support Vector Machine. A summary and conclusions are given in Section IV.

## II. Mathematics Behind Support Vector Machine

The problem of binary classification is well studied, We present an approach known as the support vector machine (SVM), which solves the binary classification task. As in regression, we have a supervised learning task, where we have a set of examples $x_n \in R^D$ along with their corresponding (binary) labels $y_n \in +1, -1$. Given a training data set consisting of example–label pairs $(x_1, y_1), ..., (x_n, y_n)$, we would like to estimate parameters of the model that will give the smallest classification error. We consider a linear model, and hide away the non-linearity in a transformation $\phi$. There are two main reasons why we chose to illustrate binary classification using SVMs. First, the SVM allows for a geometric way to think about supervised machine learning. In the papers so far, we have considered the machine learning problem in terms of probabilistic models and attacked it using maximum likelihood estimation and Bayesian inference. Here we will consider an alternative approach where we reason geometrically about the machine learning task. The SVM view of machine learning is subtly different from the maximum likelihood view. The maximum likelihood view proposes a model based on a probabilistic view of the data distribution, from which an optimization problem is derived. In contrast, the SVM view starts by designing a particular function that is to be optimized during training based on geometric intuitions. Intuitively, we imagine binary classification data, which can be separated by a hyperplane. Here, every example $x_n$ (a vector of dimension 2) is a two dimensional vector ($x_n^1$ and $x_n^2$), and the corresponding binary label $y_n$ is one of two different symbols. A hyperplane is an affine subspace of dimension $D-1$ (if the corresponding vector space is dimension D). The examples consist of two classes (there are two possible labels) that have features (the components of the vector representing the example) arranged in such a way as to allow us to separate/classify them by drawing a straight line.

### A. A. Separating Hyperplanes

Given two examples represented as vectors $x_i$ and $x_j$, one way to compute the similarity between them is using an inner product $< x_i; x_j >$. Inner products are closely related to the angle between two vectors. The value of the inner product between two vectors depends on the length

(norm) of each vector. Furthermore, inner products allow us to rigorously define geometric concepts such as orthogonality and projections. The main idea behind many classification algorithms is to represent data in $R^D$ and then partition this space, ideally in a way that examples with the same label (and no other examples) are in the same partition. In the case of binary classification, the space would be divided into two parts corresponding to the positive and negative classes, respectively. We consider a particularly convenient partition, which is to (linearly) split the space into two halves using a hyperplane. Let example $x \in R^D$ be an element of the data space. Consider a function
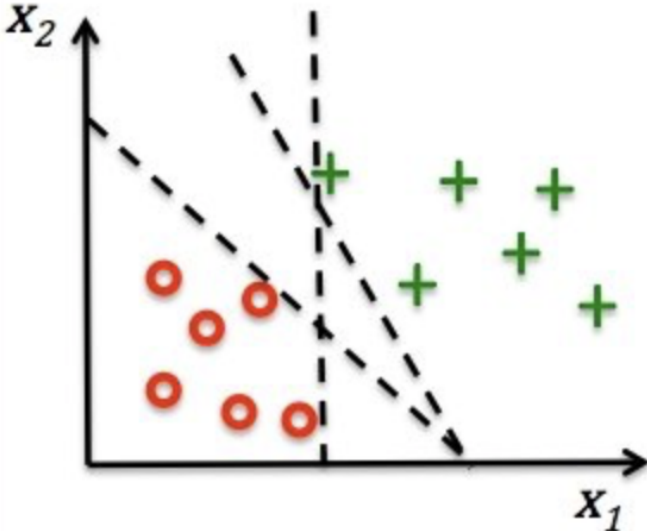
$$f : R^D \to R$$

$$x \to f(x) := < w, x > +b$$

parametrized by $w \in R^D$ and $b \in R$. Hyperplanes are affine subspaces. Therefore, we define the hyperplane that separates the two classes in our binary classification problem as

$$x \in R^D : f(x) = 0$$

When presented with a test example, we classify the example as positive or negative depending on the side of the hyperplane on which it occurs. Note that above equation not only defines a hyperplane; it additionally defines a direction. In other words, it defines the positive and negative side of the hyperplane. Therefore, to classify a test example $xtest$, we calculate the value of the function $f(xtest)$ and classify the example as $+1$ if $f(xtest) > 0$ and 1 otherwise. Thinking geometrically, the positive examples lie above the hyperplane and the negative examples below the hyperplane. When training the classifier, we want to ensure that the examples with positive labels are on the positive side of the hyperplane, i.e., $< w, x_n > +b \geq 0$ for $y_n = +1$ and the examples with negative labels are on the negative side, i.e., $< w, x_n > +b < 0$ for $y_n = 1$ These two conditions are often presented in a single equation $y_n(< w, x_n > +b) \geq 0$

*B. B. Primal Support Vector Machine*



Based on the concept of distances from points to a hyperplane, we now are in a position to discuss the support vector machine. For a dataset $(x_1, y_1), ..., (x_N, y_N)$ that is linearly separable, we have infinitely many candidate hyperplanes (refer to above figure), and therefore classifiers, that solve our classification problem without any (training) errors. To find a unique solution, one idea is to choose the separating hyperplane that maximizes the margin between the positive and negative examples. In other words, we want the positive and negative examples to be separated by a large margin (Sec. II-C). In the following, we compute the distance between an example and a hyperplane to derive the margin. The closest point on the hyperplane to a given point (example $x_n$) is obtained by the orthogonal projection.

*C. C. Concept of the Margin*

The concept of the margin is intuitively simple: It is the distance of the separating hyperplane to the closest examples in the dataset, assuming that the dataset is linearly separable. However, we need to define a scale at which to measure the distance. A potential scale is to consider the scale of the data, i.e., the raw values of $x_n$. There are problems with this, as we could change the units of measurement of $x_n$ and change the values in $x_n$, and, hence, change the distance to the hyperplane. As we will see shortly, we define the scale based on the equation of the hyperplane itself.
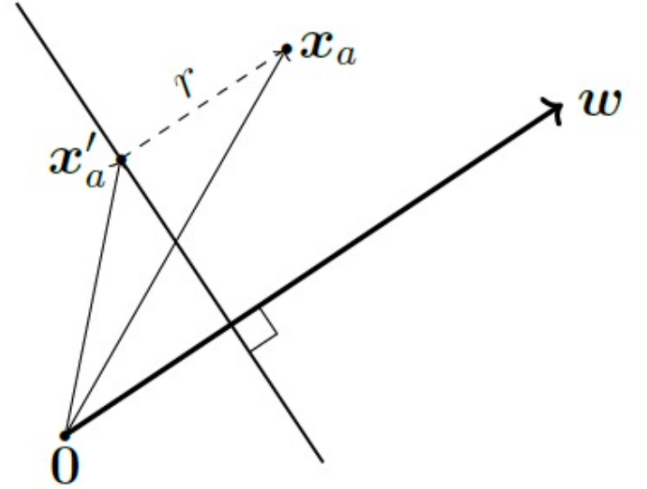


Fig. 2. Vector addition to express distance to hyperplane

Consider a hyperplane $< x_i + b >$, and an example $x_a$. Without loss of generality, we can consider the example $x_a$ to be on the positive side of the hyperplane, i.e., $< h_w; x_a > +b > 0$. We would like to compute the distance $r > 0$ of $x_a$ from the hyperplane. We do so by considering the orthogonal projection of $x_a$ onto the hyperplane, which we denote by $x_a$. Since $w$ is orthogonal to the hyperplane, we know that the distance $r$ is just a scaling of this vector $w$. If the length

of $w$ is known, then we can use this scaling factor $r$ factor to work out the absolute distance between $x_a$ and $x_a$. For convenience, we choose to use a vector of unit length (its norm is 1) and obtain this by dividing w by its norm ($\frac{w}{|w|}$), Using vector addition, we obtain

$$x_a = x_a + r\frac{w}{|w|}$$

we would like the positive examples to be further than r from the hyperplane, and the negative examples to be further than distance r (in the negative direction) from the hyperplane. Analogously to the combination of (4) and (5) into (6) we formulate this objective as

$$y_n(< w, x_n > +b) \geq r$$

Since we are interested only in the direction, we add an assumption to our model that the parameter vector w is of unit length, i.e., $||w|| = 1$, where we use the Euclidean norm $||w|| = w^T w$. This assumption also allows a more intuitive interpretation of the distance $r$ since it is the scaling factor of a vector of length 1. Collecting the three requirements into a single constrained optimization we get our objective function as $min_{w,b}\frac{1}{2}||w||^2$ subject to $y_n(< w, x_n > +b) \geq 1$

## III. THE PROBLEM

The given dataset describes a sample of pulsar candidates discovered during the survey. Pulsars are a rare type of Neutron star that emits radio waves that can be detected on Earth. They are highly scientifically interesting as probes of space-time, the interstellar medium, and states of matter. As pulsars rotate, their emission beam sweeps across the sky, producing a detectable pattern of broadband radio emission when it crosses our line of sight. Because pulsars rotate so quickly, this pattern repeats itself regularly. Thus, pulsar search entails using large radio telescopes to look for periodic radio signals. Each pulsar emits a slightly different pattern, which varies somewhat with rotation. As a result, a potential signal detection known as a "candidate" is averaged over many rotations of the pulsar, as determined by the length of observation.

### A. A. Data Description

In the absence of additional information, each candidate could describe a real pulsar. In practice, however, almost all detections are caused by radio frequency interference (RFI) and noise, making legitimate signals challenging to detect. To facilitate rapid analysis, machine learning tools are now being used to automatically label pulsar candidates. Classification systems, in particular, that treat candidate data sets as binary classification problems are becoming increasingly popular. In this case, legitimate pulsar examples constitute a minority positive class, while spurious examples constitute the majority negative class. The data set shared here includes 16,259 spurious examples caused by RFI/noise and 1,639 genuine pulsar examples. All of these examples have been reviewed by human annotators.
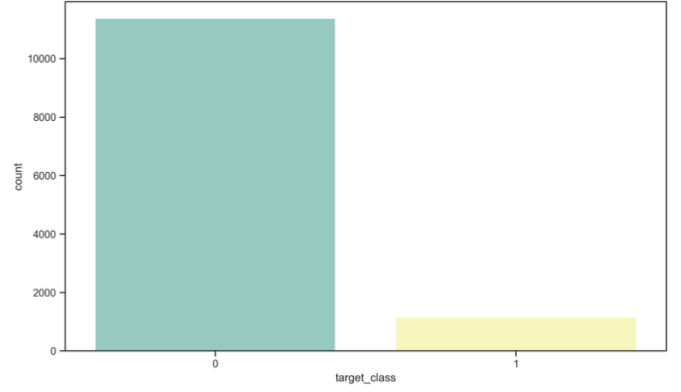


Fig. 3. target_class distribution

Eight continuous variables and one class variable are used to describe each candidate. The first four are detailed statistics derived from the integrated pulse profile (folded profile). The remaining four variables are derived in the same way from the DM-SNR curve.

The aim of the analysis is to determine whether is star is a pulsar star or not base using a generalised support vector machine classification model. Once the data is imported it is highlighted that there are missing values in the columns containing information about Excess kurtosis of the integrated profile (13.85%), Standard deviation of the DM-SNR curve (9.402%) and Skewness of the DM-SNR curve (1.8%). From Fig. 3 it is observed that most stars belong to target class 0. The number of missing values is significantly less. Therefore they had been imputed with the corresponding median values.

### B. B. Exploratory Data Analysis

There are total 12528 observations in the train set and 5370 observations in the test set, with no duplicate rows anywhere. The target is NaN in test set as it is to be predicted. All the columns are numerical in nature. The data contains positive, negative, as well as null values. Some important observations we obtain pertaining to the data are regarding the mean and distribution of each variable.
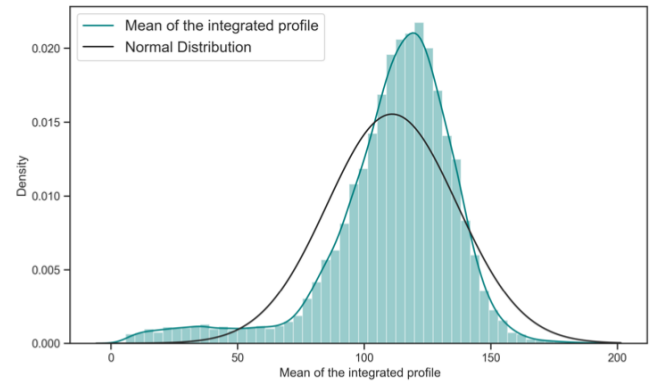


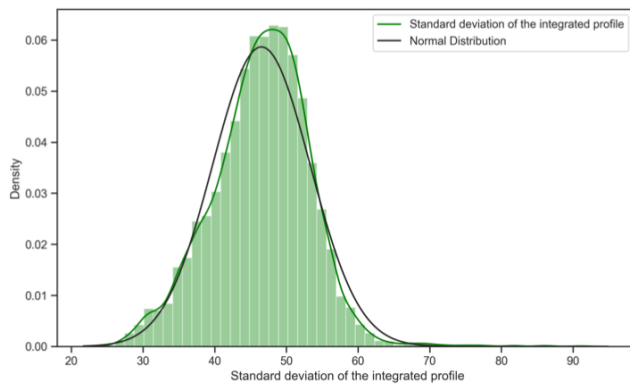Fig. 4. Distribution of *Mean of the integrated profile*

Fig. 5. Distribution of *Standard deviation of the integrated profile*



Fig. 7. Distribution of *Skewness of the integrated profile*

The distribution of Mean of the integrated profile is shown in Fig. 4. We observed that the mean of the data in the column is 111.0418 with a high standard deviation of 25.672828. The figure also tells us that the data is slightly far from a normal distribution and possesses negative skewness. The Standard deviation of the Integrated profile has a near- normal distribution, as seen in Fig. 5. The mean SD is 46.6, close to the median and looking at min and max values, it seems to be generally distributed throughout the dataset. In the Excess Kurtosis (Missing Values) majority of the data (atleast 75%) is less than mean. Hence a large head portion in this distribution. Hence the distribution of the left of mean is more tightly spread than the right. This means the integrated profile's tails are generally the same size as normal distributions.
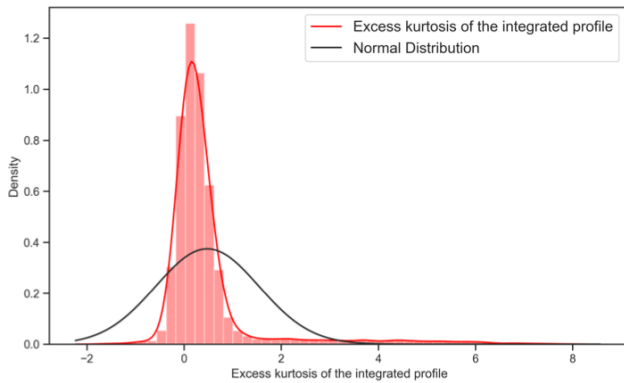
For the DM-SNR Curve the Mean is 12.674 and the standard deviation is pretty high i.e. 29.613, with with more than 75% values being below 5.6 and the max value extremely high. We can expect the mean of most curves to be on the lower side, as seen in Fig. 8.



Fig. 8. Distribution of *Mean of the DM-SNR Curve*



Fig. 6. Distribution of *Excess Kurtosis of the integrated profile*



Fig. 9. Distribution of *Standard Deviation of the DM-SNR Curve*

The majority of this data (definitely more 75%) is less than the mean. Hence a large head portion in this distribution. Therefore the distribution of the left of mean is more tightly spread than the right. Thus the integrated profile must not be very skewed. The values are pretty high, mean being 102 and the max being 1191, which is exceptionally high compared to the 75th quartile.
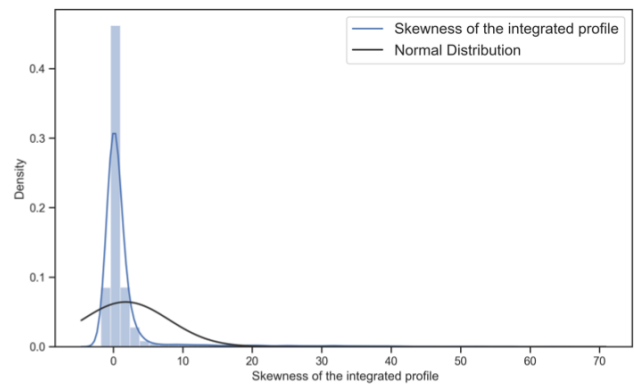
As seen in Fig. 9 the Standard deviation of DM-SNR Curve

also appears pretty skewed as maximum value is pretty high at 110.64 and 75% of all values are below 28, the median is 19, mean is 26, mean's higher magnitude could be a result of the extremely high values, but if they were omitted, then we could expect the standard deviation to be around the median value. The Excess kurtosis values are pretty high, mean being 8.23 and the max is 34, the values seem to be normally distributed with mean nearly equal to 8.32 and standard deviation of 4.53.

Following figure describes the distribution of each data variable with respect to the target class.



## C. Model Fitting

The python package `sklearn` is used for performing the task of fitting a Support Vector Machine classifier. Scikit-learn (Sklearn) is one of the most valuable and robust libraries for machine learning in Python. It provides a selection of efficient tools for machine learning and statistical modelling, including classification, regression, clustering and dimensionality reduction via a consistent interface in Python. This library, which is primarily written in Python, is built upon `NumPy`, `SciPy` and `Matplotlib`. The sklearn package provides different classes for the classification problem, including svm. Before proceeding with the modelling, all the variables are scaled using and transformed to mean 0 and standard deviation as 1. In short, we standardize the data. It is helpful for data that has negative values. It arranges the data in a standard normal distribution. Once that is done, the data is split into train and test data. Afterwards, SVM is fitted on

the training data, and the labels for test data is predicted. The mean accuracy score on the test data using basic Support Vector Machine using Radial Basis funtion Kernal is 0.9757 where as the test accuracy is 0.90774, which indicates that our model is performing decently. The majority of stars lies in the category 0. The null accuracy obtained is 0.9084, which therefore doesn't convey information about whether the classifier is doing a decent job or predicting the majority class only. The 10 fold cross-validation score for the training data is 0.9741 and for the test data it is 0.9753. The confusion matrix describes the following:

- True Positives (TP) = 3389
- True Negatives (TN) = 279
- False Positives (FP) = 26
- False Negatives (FN) = 65

The train and test classification report suggest the following

|          | precision | recall | f1-score | support |
|----------|-----------|--------|----------|---------|
| **0**    | 0.98      | 0.99   | 0.99     | 7960    |
| **1**    | 0.93      | 0.80   | 0.86     | 809     |
|          |           |        |          |         |
| **accuracy** |       |        | 0.98     | 8769    |

TABLE I
CLASSIFICATION REPORT OF THE TRAINING DATA

Since the results we obtained are satisfactory enough so we trained the data as whole and predicted the value on the given test set after performing scaling and removing the outliers.

|          | precision | recall | f1-score | support |
|----------|-----------|--------|----------|---------|
| **0**    | 0.98      | 0.99   | 0.99     | 3415    |
| **1**    | 0.91      | 0.81   | 0.86     | 344     |
|          |           |        |          |         |
| **accuracy** |       |        | 0.98     | 3759    |

TABLE II
CLASSIFICATION REPORT OF THE TEST DATA

## IV. CONCLUSIONS

To date, only over 2,000 pulsars have been detected. The task of classifying whether a star is a pulsar or not is beneficial for scientists as well as astronomers for their celestial research. The given data is used above to describe the relationship between different variables and identify the most useful ones. Before modelling, we dropped the variables like Excess kurtosis of the integrated profile, Skewness of the DM-SNR curve and Standard deviation of the DM-SNR curve due to their high correlation with other variables. The EDA and Support Vector classifier approach gives a good idea about the relationship between predictors and regressors. It explains the underlying pattern thoroughly and predicts well. Future work will focus on finding the best imputation approach and

performing hyperparameter selection using `GridSearchCV` and `RandomSearcCV`.

## REFERENCES

[1] Christopher M. Bishop, "Pattern Recognition and Machine Learning"
[2] Kevin P. Murphy, "Machine Learning, A Probabilistic Perspective"
[3] Yunus Koloğlu, Hasan Birinci, Sevde Ilgaz Kanalmaz, Burhan Özyılmaz, "A Multiple Linear Regression Approach For Estimating the Market Value of Football Players in Forward Position"