

# EMPLOYABILITY OF WEB MINING APPLICATION TO CLASSIFY BIOINFORMATICS DATASETS

SARTHAK GARKHEL

## ABSTRACT

*A web application is worked for order bioinformatics datasets. Our application gives a simple and intuitive visual interface which will be valuable for specialized and non-specialized clients. This application is mainly used for classification bioinformatics datasets, especially multi class large datasets, using sequential and parallel classification algorithms that is hopefully be widespread acceptance and adopted in both academia and business. Biological datasets are applied and classified using both serial as well as parallel support vector machine. Our proposed application has been changed altogether without any preparation introduces a general system for information pre-handling, order, and expectation. These three main tasks are applied in different datasets of different size such as Leukemia, Colon-cancer, Breast-cancer, DNA, and Protein. In the pre-processing phase, various types of data pre-processing techniques like Data Cleaning, Data Transformation, Data Reduction, and Data Discretization are used to solve incomplete and/or inconsistent problems in raw data. Then, in classification phase, a classification starts to work on pre-processing data according to different algorithms such as Serial SVM Algorithm, Parallel SVM Algorithm, Clustering, Decision Trees, Genetic Programming, and Bayesian Networks to produce a trained model based on training datasets. Finally, in the prediction phase, the trained model is used to predict the class value of a new instance in a given dataset. In order to establish an efficient and effective prediction model, we have taken into account that our prediction model must have the following criteria Accuracy, Speed, Robustness, and Scalability. The purposed application has shown much promise due to its robust classification capabilities to produce a prediction model with high accuracy ranging from 70.32 % to 97.33 %.*

## 1. INTRODUCTION

Classification considered as one of the most significant techniques that analyses data which can be apply for data classes' classification. Data classification passes through two steps. As the initial step, it constructs a model through preparing a known arrangement of information classes; that is named preparing dataset. In the subsequent advance, the scholarly inferred model is tried utilizing testing information, in which tests are haphazardly and vary from the preparation tests, for assessing the precision of this model by means of contrasting the class marks that are anticipated by the learned model with the realized class names. The precision of a model estimated by the level of test tests that are accurately arranged by the model. As a rule, the took in model got from the first1. When all is said in done, the took in model got from the initial step can be choice trees, scientific formulae or grouping rules. There are numerous strategies for information characterization, for example, choice tree enlistment, Bayesian grouping, Rule-Based arrangement, Neural Networks, and Support Vector Machine (SVM). However, SVM is highly used and the most famous one. SVM is very useful and effective among all machine learning algorithms and depends on theory of statistical learning. In general, there are two kinds of machine learning techniques, supervised and

non-supervised, used for pattern recognition<sup>2</sup>. In term of SVM, it is considered as a supervised machine learning technique where class label must be known in advance<sup>3</sup>.

The main function of SVM algorithm is to divide two point classes of a learning data set with a surface. With this classification algorithm, the margin can be increased between two points classes<sup>2</sup>. Assume that training data set D has the form

$$D = \{(X_i, Y_i), i = 1, 2, \dots, n, X_i \in \mathbb{R}^M, Y_i \in \{-1, +1\}\} \quad (1-1)$$

where,  $X_i$  is data sample and  $Y_i$  is label. The objective is to find a function  $F: \mathbb{R}^M \rightarrow \mathbb{R}$  such that, for any sample  $X_i$  there are trained parameters  $W_i$  and  $b$  where,

$$F(X_i) = X_i \cdot W_i + b = \sum_{i=1}^m X_i W_i + b \quad (1-2)$$

Then, for any sample  $(X, Y)$ , the class assigned by the model is

$$Y = \text{Sign } F(X) = \begin{cases} -1 & \text{if } F(X) \leq 0 \\ +1 & \text{if } F(X) > 0 \end{cases} \quad (1-3)$$

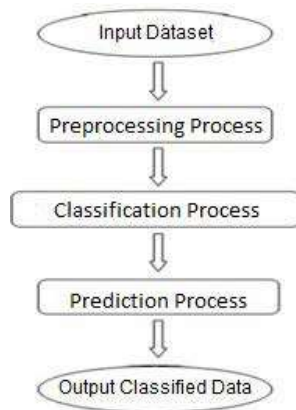
time for task booking and consolidating the outcomes. The proportion of the parallel execution with the relative speedup ( $S$ ) is characterized as the proportion of time expended in taking care of an issue on one processor to the time expected to take care of a similar issue on a parallel PC machine with  $P$  indistinguishable handling component. Effectiveness ( $E$ ) is likewise a significant method to assess the parallel usage, named as the proportion of speedup to the processors number<sup>3</sup>. SMC-PBC-SVM and PMC-PBC-SVM are considered as a lot of proficient parallel calculations. These two calculations can run well to join Parallel Binary Classes with Serial Multi-Class SVM for classification<sup>4</sup>, while PMC-PBC-SVM consolidates Parallel Multi-Class Support Vector Machine with Parallel Binary Classes for classification<sup>5</sup>. The central thought of two calculations is to isolate a few processors in to two separate subsets. One is utilized for a multi-class case and the other is for the parallel class case. The multi-class case bunch in SMC-PBC-SVM contains a single processor used to fathom multi-case in sequential though the twofold class bunch contains the rest of the processors, used to illuminate paired case in parallel<sup>4</sup>. While in PMC-PBCSVM calculation, a lot of processors are partitioned as a lattice with lines. Each line can settle differed paired case in parallel that implies multi-case is at long last sifted through in parallel<sup>5</sup>. Electronic applications are those applications that actualize on the web. These kinds of utilizations need an internet browser. These applications are made in the program upheld programming language, like JSP, JavaScript, HTML, CSS and so on. WEKA is a case of these sorts of utilizations for classification<sup>6</sup>. Web applications generally utilize a mix of server-side content (JSP, PHP, and so on.) and customer side content (HTML, JavaScript, and so on.) to build up the application. The customer side content arrangements with the introduction of the data while the server-side content arrangements with all the hard stuff like putting away and recovering data. The World Wide Web is fast budding as a

significant medium for business as well as for the spreading of information related to a wide variety of topic (e.g., industry and administration). According to most prediction, the bulk of individual information will be obtainable on the Web. These vast amount of data elevate a majestic confront, namely, how to turn the Web into a more practical information service<sup>7</sup>. These days there is a huge quantity of applications and services that are on hand from side to side Internet as they are looking for, chat, sale, etc., yet a large amount of that information is not helpful for a lot of people, but in the field of Data Mining, all the information available in the Internet represent a work prospect and it is possible to do many study on the foundation of these with detailed purpose. Data Mining and Knowledge Discovery Data, KDD, are smart tools for data analysis where the speedy distribution of these technologies calls for an urgent test of their communal crash. A summary of KDD and Data mining technologies have been shown. The two Terms Data Mining and KDD are used to explain the 'non-trivial extraction of understood, before indefinite and potentially helpful information from data<sup>8</sup>. KDD is an idea that describes the procedure of search on huge volume of data for pattern that can be measured knowledge about the data<sup>9</sup>. The most famous division of knowledge discovery is data mining.

A model has been proposed to optimize a big data by authors in<sup>14</sup>; they used the coordinate system to improve the result Pre-processing techniques have been used to obtain highly accurate Classification results due to real data often are not ready for processing directly because they are may not completed (missing field values, missing some important attributes, or having aggregate data only), noisy (having outliers or mistakes) or not consistent (having contrarities in codes or names<sup>15</sup>. Major techniques used in pre-processing for data are data cleaning in which many techniques can be used such as replacing lacked values with real values, smoothing noisy data, determining or dropping outliers, and resolving contrarities<sup>13</sup>, data integration means integrating of several databases<sup>22</sup>, data transformation means aggregating and normalizing<sup>16</sup>, by which minimizing is implemented for representing in volume but also give us the same or analogous analytical results<sup>17</sup>, and Data discretization belongs to data reduction but especially for numerical data<sup>18</sup>.

## 2. WEB-BASED PROPOSED APPLICATION

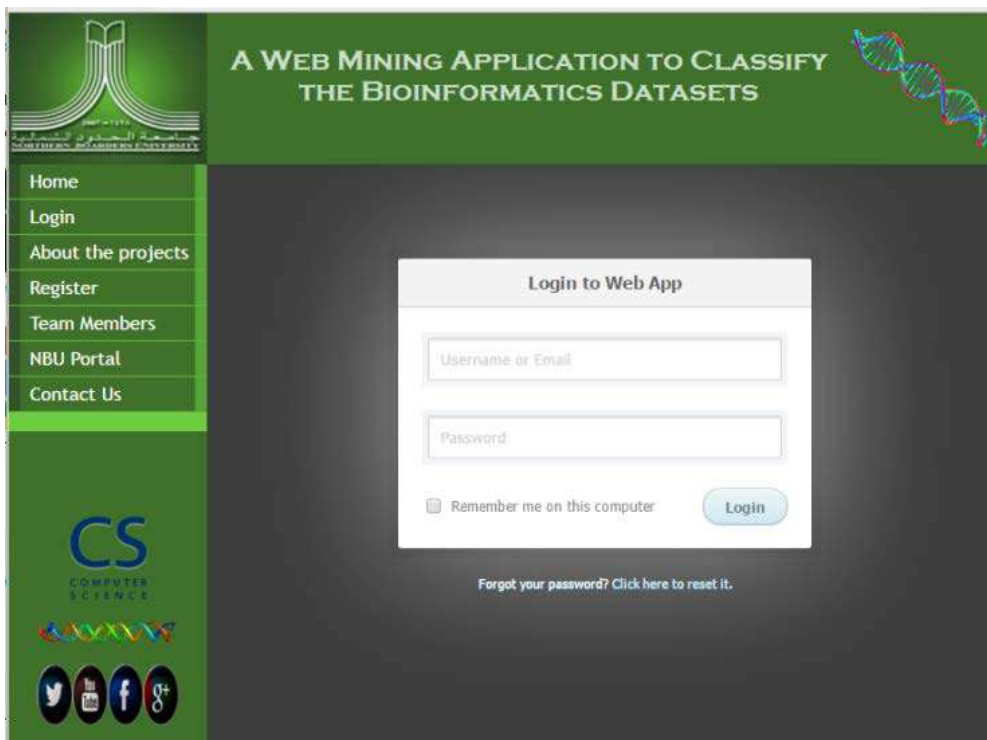
There are many fields where classification can be used such medical, marketing, trading and science purposes, etc. Our website is a web mining application that allows any user to run classification even sequential or parallel for any dataset in order to reduce the time of analysing and increase prediction accuracy without needing to install a complex Software on his/her computer that may require a higher CPU speed and a larger Memory size. The proposed model consists of three main tasks:



**Figure 1.** Proposed model of web mining.

Pre-processing, Classification and predication as shown in Figure 1.

JSP technology is used to implement this model as a web application, which let a user to create his/her own account if he/she has no account as shown in Figure 2 and to use this web application many times later via the account that has been created. The main tasks will be explained in details in the following:



**Figure 2.** Login screen.

### 3. DATA PRE-PROCESSING TECHNIQUES

Pre-processing of information is considered as a procedure of information mining that includes changing crude information into a reasonable organization. Certifiable information can be deficient, conflicting and/or come up short on specific practices or patterns. It is probably going to contain

numerous blunders. Information pre-preparing has been demonstrated to determine such issues. Cutting edge databases are helpless to uproarious, absent and conflicting information. This is a direct result of their gigantic size (frequently a few gigabytes or then again more) and in light of the fact that they typically originate from various, heterogeneous sources. Low-quality information will bring about low-quality information mining results. In this way, Pre-preparing is one of the principal steps that improve the precision of the outcomes. In the event that we have great pre-handling, at that point we will accomplish right and precise outcomes. We will have more control over methods for information treatment. This is normally done by entering the information into an exhibit which enables us to alter it easily as we wish.

Data mining is the extraction of specific information from large data sources. The KDD abbreviation refers to data mining which stands for Knowledge Discovery in Databases, Knowledge Extraction and Data/pattern analysis refers also to data mining. KDD has 7 stages as part of an iterative sequence, which are used with specific methods<sup>19</sup>. The first 4 stages are: Cleaning, Integration, Transformation (including reduction and Discretization methods) and the Selection of data. These are part of data pre-processing. After these stages, there are 3 more stages, which are different types of mining processes. These involve Data mining and Pattern evaluation as well as presentation of knowledge<sup>20</sup>.

In the world that we live in, various types of data including similar types of data are grouped together in one place with different data sources. A lot of data is lost or irregular due to missing data, noise, inconsistent or even outliers. Such data does not produce high quality information. If such low quality data is produced in terms of classification, data analysis, pattern reorganization and decision management then that means that data mining has not been used to produce the optimum solution. In order to produce quality data by removing irregularities,

data pre-processing techniques must be used for data mining tasks. Many data pre-processing techniques are used as part of data pre-processing methods, in order to remove specific irregularities. Data cleaning, Data Integration and Transformation, Data reduction and Discretization are all basic data pre-processing methods.

In this section, the ways in which data pre-processing works will be analysed and described. Data pre-processing has various types of methods and technique in order to function. As shown in Figure 3, a user can choose a suitable method for pre-processing his data.



**Figure 3.** Preprocessing methods screen.

### 3.1 Data Cleaning

It is the first stage of data pre-processing methods. One of the main problems with building data warehousing and mining is that real world data is corrupt due to noise and missing values in tuples and is inconsistent. Applying mining techniques to such data will result is unreliable and poor output. Data needs to be “cleaned” before data mining. This is done by reducing noisy “interference” data, filling in missing values, correcting inconsistencies and identifying or removing outliers. Different types of data cleaning techniques are described below<sup>21,22</sup>. In order to tackle Missing Values in datasets, there are many techniques. First, is to remove the blank tuples. This is not an efficient technique, for example, any tuples that have blank attributes are ignored. Second, manually filling in missing value, for example, filling in blank tuples with assuming what the most suitable values for their attributes, but manually inputting values is not efficient for large data sets and is considered time-consuming. Third, a global constant can be used to resolve the problem of missing values, where all missing values are replaced with the same or global constant. Fourth technique, calculating the attribute means to replace all the missing attribute values with the mean (for numerical data). This technique identifies the mean for unfilled attributes and fill use this mean value for unfilled tuples, for example, identify the mean for specific numerical attributes and use this for every blank tuple that has specifies attributes. Finally, the most probable values are used to fill in the missing value where the most recurring values are filled. Noisy Data can be tackled by different methods. First, Binning (used for numerical data) where attribute values are distributed into numerous bins. Each value in the bin is then replaced by the mean value for the bin. This process is known as “smoothing by bin means. Second, Regression (used for numerical data); this method involves finding a line of best fit for two attributes, where one attribute is used to predict the other. Finally, Clustering contains similar values in forms of groups or clusters and some of these values fall outside these groups or clusters, in order to know and show the outliers, the values are clustered into groups, where the values that founded outside the groups are considered to be outliers.

### 3.1.1 Data Transformation

Change of various sorts of information starting with one arrangement then onto the next is the biggest issue that should be conquered while developing things, for example, information mining, information warehousing, and the World Wide Web. This is because real world data is in different formats and language. If mining or analysis technique are applied to the data in data warehouses or WWW, they will take require more time and occupy more memory. Decisions will also be suspended and the quality will be less. Data would then need to be transformed from one format to another by using aggregation, generalization and normalization techniques<sup>23</sup>. Various transformation techniques can be used. First, Aggregation (used for numerical data): where data is calculated to summarize a specific attribute in this technique. This technique is useful for constructing data cubes for analysing multiple granularities. Second, Generalization, also popularly referred to as concept hierarchy technique, it is higher level data and it is used to replace lower level data. This technique is useful for categorical and numerical attributes, Finally, Normalization: this involves scaling the data range to include a specific range such as [0.0, 1.0]. This technique is useful for ANN classification algorithms, Normalization: this involves scaling the data range to include a specific range such as [0.0, 1.0]. This technique is useful for ANN classification algorithms. There are many methods for data normalization. Min-Max normalization, z-score normalization, and normalization by decimal scaling are some important types of them.

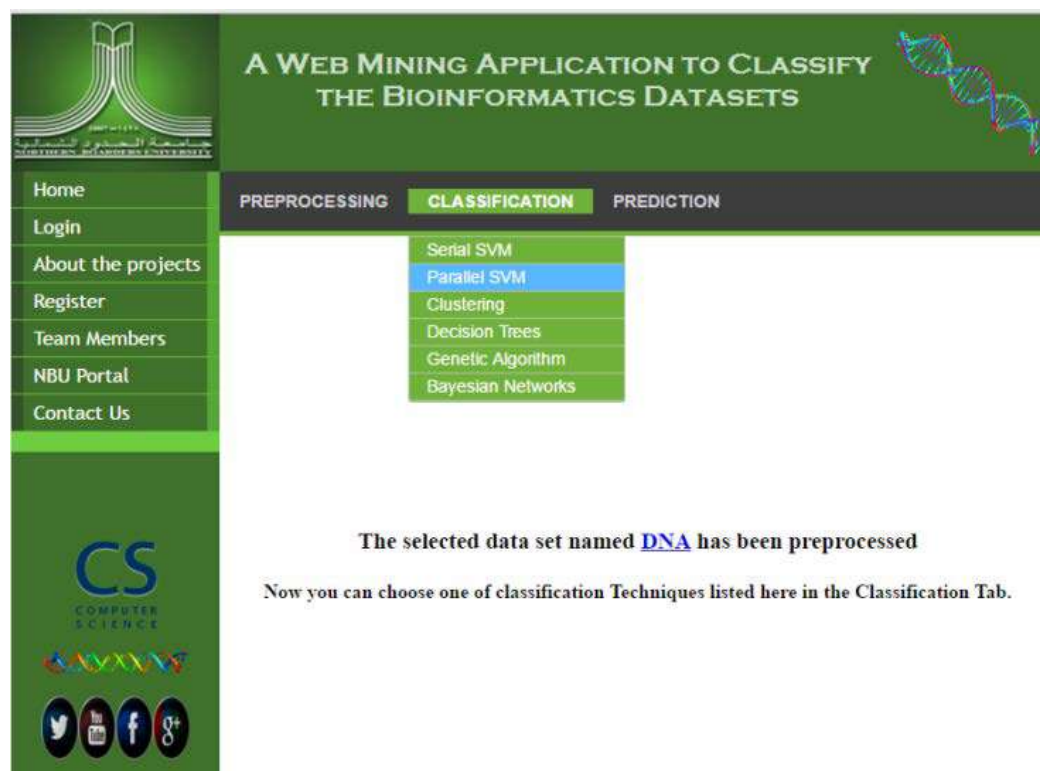
### 3.1.2 Data Reduction

Databases and data warehouses may store large amount of data (such as 1 TB). Therefore, nowadays the daily increase in the magnitude of the data poses a problem. Data that has been mined and even simple data are store in various places such as databases, data warehouses and the WWW. The mining techniques applied to these types of data will consume too long time in order to process. Their complexity would also form part of the decision or result. Reduce data volumes are obtained via aggregation of data cube, reduction of dimensionality, Binning, Binning and Sampling technique. Data Reduction can be applied to number of attributes, attribute values, and number of tuples<sup>24</sup>. Many methods are used to reduce number of attributes. First, Data cube aggregation, this method is useful for constructing data cubes. Data cubes store multidimensional aggregated information. Each cell carries an aggregate data value, corresponding to multi-dimensional spaces. Second, Dimensionality reduction, data encoding are applied to obtain a reduced representation of the original data. Finally, Attribute subset selection: problems such as irrelevant, weakly relevant and redundant attributes are detected with this and removed. This results in reduced data set sizes and is achieved by using greedy algorithms, decision tree induction or Metaheuristics Algorithms such as Ant colony and Genetic Algorithms. Attribute selection is to get a smallest subset of attributes, where results satisfy probability distribution of data classes which implies to be as close as possible to the original distribution obtained through all the attributes<sup>19</sup>.

### 3.2 Classification Techniques

classification is one of the most significant methods of information mining. It is valuable in various fields of research like forecast, data recovery, web searches and that's just the beginning. Most present clustering strategies are sequential, which isn't compelling in enormous datasets. Along

these lines, there is a requirement for Parallel Classification, which orders enormous informational index by partitioning it among numerous processors so as to speed up arrangement. We can group to our information in such a large number of ways, as Matrix technique, Decision Tree, SVM, k-neighbours and so on. The fundamental burden of SVM is that to require the enormous memory for the calculation. Our framework can arrange the enormous information online on an extremely ostensible necessity memory. In the wake of getting pre-handled information from the pre-preparing step, parallel characterization begins to work as indicated by PMC-PBC-SVM calculation. An enlisted client can play out the grouping step as appeared in Figure 4. After applying classification step, we created a trained model which is based on the training database. That is competitively easy problem if it has only a few important parameters, like, if we have two parameters one is simply create a scatter-plot of the quality values and the other can establish graphically how to divide the plane into uniform regions while the objects are of the same classes. On the other hand, the classification problem is called hard, when there are many parameters to consider. It is not only to visualize the dimensional space, but there are so many different groups of parameters that techniques based on exhaustive searches of the parameter space quickly become computationally infeasible. Practically for classification always involve a best approach proposed to find a good enough solution to the optimization problem. There are so many classification methods, but we are discussing here few of them.



**Figure 4.** Classification process screen.

### 3.2.1 Serial SVM Algorithm

This algorithm works on the contender Support Vector set. It starts the set with the nearby couple of points from differing classes similar to the Direct SVM algorithm. As soon as the algorithm finds



a violating point in the dataset it greedily adds it to the candidate set. It may so happen that adding up the violate point as a Support Vector may be prohibited by other contender Support Vectors by now present in the set. We simply clip away all such points from the contender set. To make sure that the KKT environment are fulfilled we make continual passes through the dataset unless no any violators can be found. We employ the quadratic penalty formulation to make sure the linear partition of the data points in the kernel portion<sup>28</sup>.

### 3.2.2 Parallel SVM Algorithm

Another method for grouping calculation is presented as most of the results are simply promising, so analysts need to some additional work to make it increasingly helpful. To build up a skilled and supportive parallel calculation for an arrangement which is named PMC-PBC-SVM,

the calculation is a consolidated double arrangement issue, which is settled in parallel, with a multi-class characterization issue, which is fathomed in parallel. The primary goal of this calculation is to separate a lot of processors into two subsets. The main arrangement of processors is utilized to take care of the multi-class issues in parallel. While other arrangements of processors are answerable for taking care of the double class issues in parallel as well. By the PMC-PBC-SVM calculation, the analyzation of a gathering of processors resembles as a network where each line is utilized to chip away at an alternate paired case in parallel.

Consequently, the multi-class case hopes to tackled in parallel. It empowers us to do two-phase parallelisms. This calculation, named PMC-PBC-SVM, is helping us to join a parallel Multi-Class Support Vector Machine with Parallel Binary Class for an order. It surely turns out to be progressively proficient and viable by the accompanying algorithm<sup>5</sup>. The SVM calculation investigations a lot of marked information tests to have the option to arrange new information tests. The parallel SVM calculation has been created to upgrade the adaptability. The SVM examinations a gathering of marked examples so as to group new example information. The primary thought of PSVM is to diminish computational expenses by applying parallel ICF which burdens preparing informational collection onto parallel machines and runs factorization for every one of them on these machines. When PICF has stacked in preparing information spread on  $m$  machines and diminished the size of the part grid through factorization, IPM can be explained on parallel machines concurrently<sup>29</sup>.

### 3.2.3 Clustering

Clustering allows a user to make groups of data to determine patterns from the data. Clustering has its advantages when the data set is defined and a general pattern needs to be determined from the data. We can create a specific number of groups, depending on our business needs. One defining benefit of clustering over classification is that every attribute in the data set will be used to analyse the data<sup>30</sup>.

### 3.2.4 Decision Trees

A decision tree is a tool that uses chart like a tree or form of decisions and their probable consequences, including possibility of event outcome, resource expenses, and value<sup>31</sup>.

### 3.2.5 Genetic Programming

GP, which stands for Genetic Programming (GP) has been used recently in order to resolve classification problems of data mining. The widely use of GP refers to the fact that prediction rules are so naturally represented in Genetic Programming and refers to the good results produced by GP with global search problems like classification. The search space of classification has many peaks which make local search algorithms are run badly<sup>32</sup>.

### 3.2.6 Bayesian Networks

By the help of Bayesian, we can create the graph or model of possible relations between the different set of features. The Bayesian network arrangement  $S$  is a Directed Acyclic Graph (DAG) and the related nodes in  $S$  are in one-tone connection with the different features. The related arcs show informal influences between the features while the lack of probable arcs in  $S$  encodes provisional independencies<sup>33</sup>.

### 3.3 Prediction

In the prediction phase, the model which is created in classification phase is used to predict out coming results<sup>34</sup>. Predictive modelling can be applied to any kind of unidentified event, in spite of when it occurred. In a lot of cases the model is chosen on the opening of detection premise to attempt to estimate the possibility of an outcome given a set quantity of input data, like, given any email formative how likely that it is spam. Models can use one or many classifiers in trying to find out the possibility of a group of data belong to another set, say it spam. Depending on some useful limitations, predictive modelling is identical with, or mostly overlap with the ground of machine learning, as it is more usually referred to in educational or research and development context<sup>35</sup>.



**Figure 5.** Prediction algorithms screen.

The purpose of this paper is to produce a roughness prediction model of bioinformatics datasets with high precision and speed which is established by parallel classification algorithms. The prediction task is a supervised learning task where the data are used directly to predict the class value of a new instance in a given dataset. In order to establish an efficient and effective prediction model, we have taken into account that our prediction model must have the following criteria: First, Accuracy, this term means that how good the predictor can estimate the assessment of predicted feature for a new data. Second, Speed, this term tells us to the computational value in developing and using the predictor. Third, Robustness, this term refers to the skill of predictor to make accurate predictions from given strident data. Forth, Scalability, this term means that it is the skill to build the predictor proficiently. Finally, Interpretability, this term means the understand ability of the predictor i.e. to what amount the predictor understands. Prediction model predicts categorical and incessantly esteemed function. For instance, a classification form can be built to classify applications of bank to detect each one if it is secure or non-secure. Based on the prepared model which made from order step, obscure information can be anticipated utilizing the third errand of our framework and afterward yield the characterized information into yield document to the client. Many prediction algorithms have been applied to our data such as CHAID, C&R, QUEST trees, neural network, Bayesian, logistic regression and SVM. You can choose a proper algorithm form drop down menu as shown in the Figure 5.

Table 2. Comparing SVM and PSVM technique for the same dataset

Dataset Name	SMC-SBC-SVM	SMC-PBC-SVM	Accuracy
	NP=1	NP=36	
Leukemia	0.6100	0.1659	96.77
Colon-cancer	0.2400	0.083	90.32
Breast-cancer	0.6100	0.7457	97.22
DNA	2110.0000	647.2050	95.70
Protein	2427.8400	313.0594	70.32
Iris	2.0000	647.2050	97.33

#### 4. RESULT ANALYSIS

We added the case study to support our result. In this study we took the biological datasets and classify them using both serial and parallel support vector machine. The past decade has seen an unstable growth in some different field of medical and computer science research. Biological data mining has become a significant area of a new research field says bioinformatics. As the area of biological data mining is active, rich, and broad, it is not possible to cover such an important and successful theme in one subsection. The human genome is expected to include around 20,000 to 25,000 genes. Genomics is the study of genome sequences. The classification of DNA or amino acid sequence patterns that play roles in different biological functions, genetic diseases, and evolution is challenging. This requires a huge transaction of research in computational algorithms, statistics, mathematical programming, data mining, machine learning, information retrieval, and other disciplines to build up efficient genomic and proteomic data analysis tools. Also, in medical application such as Radiation Pneumonitis, where lung cancer patients who receive radiotherapy as part of their treatment are at risk radiation-induced lung injury, new methods are needed to guide physicians to prescribe a targeted dosage to patients at high risk of RP which is a potentially fatal side effect to treatment. In this work, several widely classification algorithms in the machine learning area are used to differentiate between unlike set if risks of RP. Additionally, for diabetes patients, an order strategy is made arrangements for similar patient's dependent on Support Vector Machines (SVM) classifiers that put the analysed prostate diabetes patients into the set if dangers, before performing radical prostatectomy, as indicated by their therapeutic parameters.

We are here comparing the two techniques for the same data and giving the accuracy data in Table 2. In this work, we take some random data for the two types of process when we did 1 process and when we execute 36 processes. We follow here two different algorithms in the second column it is the SSVM and in the third column PSVM. After that we gave the accuracy measurement by both of them which have been shown in our fourth column in Table 2.

## **5. CONCLUSION**

The outcomes accomplished subsequent to applying our proposed application with various informational collections from various sizes with both specialized and non-specialized clients were sure through utilizing a simple and intelligent interface of our application as appeared in the above Figures. This application was executed on the web utilizing JSP and MySQL, so it is the unadulterated online application for grouping the multiclass enormous datasets. In future work, the proposed application can be adjusted by including perception of processors execution through calculations to get progressively exact outcomes.