

## CROP YIELD PREDICTION USING MACHINE LEARNING TECHNIQUES (SMART FARMING)

<sup>1</sup>Ch. KODANDARAMU, Associate Professor

<sup>2</sup>KANAPAREDDY JITENDRAKUMARI

<sup>3</sup>KORADA RAMBABU

<sup>4</sup>KUBIREDDY SASIDHAR

<sup>5</sup>GANAGALLA RAMESH

<sup>1,2,3,4,5</sup>Miracle Educational Society Group of Institutions, Vijayanagaram, Ap, India

### ABSTRACT

As a coastal state, Tamil Nadu faces uncertainty in agriculture which decreases its production. With more population and area, more productivity should be achieved but it cannot be reached. Farmers have words-of-mouth in past decades but now it cannot be used due to climatic factors. Agricultural factors and parameters make the data to get insights about the Agri-facts. Growth of IT world drives some highlights in Agriculture Sciences to help farmers with good agricultural information. Intelligence of applying modern technological methods in the field of agriculture is desirable in this current scenario. Machine Learning Techniques develops a well-defined model with the data and helps us to attain predictions. Agricultural issues like crop prediction, rotation, water requirement, fertilizer requirement and protection can be solved. Due to the variable climatic factors of the environment, there is a necessity to have an efficient technique to facilitate the crop cultivation and to lend a hand to the farmers in their production and management. This may help upcoming agriculturalists to have a better agriculture. System of recommendations can be provided to a farmer to help them in crop cultivation with the help of data mining. To implement such an approach, crops are recommended based on its climatic factors and quantity. Data Analytics paves a way to evolve useful extraction from agricultural database. Crop Dataset has been analyzed and recommendation of crops is done based on productivity and season.

## 1. INTRODUCTION

Tamil Nadu being 7th largest area in India has 6th largest population. It is the leading producer of agriculture products. Agriculture is the main occupation of Tamil Nadu people. Agriculture has a sound tone in this competitive world. Cauvery is the main source of water. Cauvery delta regions are called as rice bowl of Tamil Nadu. Rice is the major crop grown in Tamil Nadu. Other crops like Paddy, Sugarcane, Cotton, Coconut and groundnut are grown. Bio-fertilizers are produced efficiently. Many areas Farming acts as major source of occupation.

Agriculture makes a dramatic impact in the economy of a country. Due to the change of natural factors, Agriculture farming is degrading now-a-days. Agriculture directly depends on the environmental factors such as sunlight humidity, soil type, rainfall, Maximum and Minimum Temperature, climate, fertilizers, pesticides etc.

Knowledge of proper harvesting of crops is in need to bloom in Agriculture. India has

seasons of Winter which occurs from December to March. Summer season from April to June. Monsoon or rainy season lasting from July to September and 4. Post-monsoon or autumn season occurring from October to November

Due to the diversity of season and rainfall, assessment of suitable crops to cultivate is necessary. Farmers face major problems such as crop management, expected crop yield and productive yield from the crops. Farmers or cultivators need proper assistance regarding crop cultivation as now-a-days many fresh youngsters are interested in agriculture.

Impact of IT sector in assessing real world problem is moving at a faster rate. Data is increasing day by day in field of agriculture. With the advancement in Internet of Things, there are ways to grasp huge data in field of Agriculture. There is a need of a system to have obvious analyzes of data of agriculture and extract or use useful information from the spreading data. To get insights from data, it has to be learnt.

## 2. LITERATURE SURVEY

[1] Shreya S. Bhanose, Kalyani A. Bogawar (2016) “Crop And Yield Prediction Model”, *International Journal of Advance Scientific Research and Engineering Trends*, Volume 1, Issue 1, April 2016

An agricultural sector necessitate for well defined and systematic approach for predicting the crops with its yield and supporting farmers to take correct decisions to enhance quality of farming. The complexity of predicting the best crops is high due to unavailability of crop knowledge-base. Crop prediction is an efficient approach for better quality farming and increase revenue. Use of data clustering algorithm is an efficient approach in field of data mining to extract useful information and give prediction. Various approaches have been implemented so far are worked either for crop prediction. Crop prediction model aiding farmers to take correct decision. This indeed helps in improving quality of farming and generate better revenue for farmers. Traditional clustering algorithms such as k-Means, improved rough k-Means and means++ makes the tasks complicated due

to random selection of initial cluster center and decision of number of clusters. Modified K-Means algorithm is thereby used to improve the accuracy of a system as it achieves the high quality clusters due to initial cluster centric selection.

[2] Tripathy, A. K., et al. (2011) *Data mining and wireless sensor network for agriculture pest/disease predictions.* *Information and Communication Technologies (WICT), 2011 World Congress on. IEEE.*

Data driven precision agriculture aspects, particularly the pest/disease management, require a dynamic crop-weather data. An experiment was conducted in a semi-arid region to understand the crop-weather-pest/disease relations using wireless sensory and field-level surveillance data on closely related and interdependent pest (Thrips) - disease (Bud Necrosis) dynamics of groundnut crop. Data mining techniques were used to turn the data into useful information/knowledge/relation/trends and correlation of crop-weather-pest/disease continuum. These dynamics obtained from the data mining techniques and trained through mathematical models were validated

with corresponding surveillance data. Results obtained from 2009 & 2010 kharif seasons (monsoon) and 2009–10 & 2010–11 rabi seasons (post monsoon) data could be used to develop a real to near real-time decision support system for pest/disease predictions

**[3] Ramesh Babu Palepu (2017) ” An Analysis of Agricultural Soils by using Data Mining Techniques”, International Journal of Engineering Science and Computing, Volume 7 Issue No. 10 October.**

Data mining is an approach through which in an synchronized manner we can find a workable solution that will be beneficial to increase the growth. The Farmers in agriculture sectors face a lot of issues and difficulties due to the improper understanding and implementation of the activities to enhance their growth and productivity. A large amount of data is available for analyses and scrutiny, however those related to agriculture sector is in a small quantity. Hence segregation and processing of the same from the sources has to be done with proper methodology. Places having multiple grain growth and

different soil structure makes it complex to have a perfect estimation of the crops yield both in quantity and quality. Creating a close link between the customer expectation and the producing capabilities of the agriculture sector can be win-win situation at both ends, this can be achieved with capturing data segment wise and in a structured manner. Thus, the customer will be able to fulfil his requirement as per his wish, rather than being satisfied by what is being offered to him. The application of such techniques enables us to predict and make analysis of various problems and helps farmers to make difficult farming decisions based on the conditions, soil fertility, crop duration, disease and other important factors that can result in poor yield production

**[4] Rajeswari and K. Arunesh (2016) “Analysing Soil Data using Data Mining Classification Techniques”, Indian Journal of Science and Technology, Volume 9, May.**

Soil is an essential key factor of agriculture. The objective of the work is to predict soil type using data mining classification techniques. Methods/Analysis: Soil type is predicted using data mining classification

techniques such as JRip, J48 and Naive Bayes. These classifier algorithms are applied to extract the knowledge from soil data and two types of soil are considered such as Red and Black. Findings: In this paper, Data Mining and agricultural Data Mining are summarized. The JRip model can produce more reliable results of this data and the Kappa Statistics in the forecast were increased. Application/Improvement: For solving the issues in Big Data, efficient methods can be created that utilize Data Mining to enhance the exactness of classification of huge soil data sets.

**[5] A.Swarupa Rani (2017), "The Impact of Data Analytics in Crop Management based on Weather Conditions", International Journal of Engineering Technology Science and Research, Volume 4, Issue 5, May.**

Many states face uncertainty in agriculture which decreases its production. With more population and area, more productivity should be achieved but it cannot be reached. Agricultural factors and parameters make the data to get insights about the Agri-facts. Growth of IT world drives some highlights in Agriculture Sciences to help farmers with

good agricultural information. The common difficulty present among the Indian farmers is they don't opt for the proper crop based on their soil necessities. Because of this productivity is affected. This provides a farmer with sort of options of crops which will be cultivated. Agricultural issues like crop prediction, rotation, water requirement, fertilizer requirement and protection can be solved. To implement such an approach, crops are recommended based on its climatic factors and quantity. Data Analytics paves a way to evolve useful extraction from agricultural database. Crop Dataset has been analyzed and recommendation of crops is done based on productivity and season.

**[6] Pritam Bose, Nikola K. Kasabov (2016), "Spiking Neural Networks for Crop Yield Estimation Based on Spatiotemporal Analysis of Image Time Series", IEEE Transactions On Geoscience And Remote Sensing.**

This paper presents spiking neural networks (SNNs) for remote sensing spatiotemporal analysis of image time series, which make use of the highly parallel and low-power-consuming neuromorphic hardware platforms possible. This paper

illustrates this concept with the introduction of the first SNN computational model for crop yield estimation from normalized difference vegetation index image timeseries. It presents the development and testing of a methodological framework which utilizes the spatial accumulation of time series of Moderate Resolution Imaging Spectroradiometer 250-m resolution data and historical crop yield data to train an SNN to make timely prediction of crop yield. The research work also includes an analysis on the optimum number of features needed to optimize the results from our experimental data set. The proposed approach was applied to estimate the winter wheat (*Triticum aestivum* L.) yield in Shandong province, one of the main winter-wheat-growing regions of China. Our method was able to predict the yield around six weeks before harvest with a very high accuracy. Our methodology provided an average accuracy of 95.64%, with an average error of prediction of 0.236 t/ha and correlation coefficient of 0.801 based on a nine-feature model.

[7] Priyanka P.Chandak (2017),” Smart Farming System Using Data Mining”,

**International Journal of Applied Engineering Research, Volume 12, Number 11**

Crop recommendation system or prediction system is the art of predicting crop yields to improve the production and production before the harvest actually takes place, it takes typically a couple of months in advance. Crop prediction depends on the computer programs that describe the plant-environment and the soil features interactions in quantitative terms. The soil testing will start with the collections of a soil sample from the field. The first basic principles of the soil testing is that a field can be sampled in such a way that by getting a chemical analysis of the soil sample and also majorly depend on temperature and rainfall will accurately reflect the field's true nutrient status on a particular area to help out farmers to improve the production.

### 3. EXISTING SYSTEM

Agriculture makes a dramatic impact in the economy of a country. Due to the change of natural factors, Agriculture farming is degrading now-a-days. Agriculture directly depends on the environmental factors such

assunlight humidity, soil type, rainfall, Maxim um and Minim um Temperature, climate, fertilizers, pesticides etc. Knowledge of proper harvesting o f crops is in need to bloom in Agriculture. India has seasons of Winter which occurs from December to March. Summer season from April to June. Monsoon or rainy season lasting from July to September and. Post-monsoon or autumn season occurring from October to November.

Due to the diversity of season and rainfall, assessment of suitable crops to cultivate is necessary. Farmers facemajor problem s such as crop management, expected crop yield and productive yield from the crops. Farers orcultivators need proper assistant regarding crop cultivation as now-a-days many fresh youngsters are interestedin agriculture.

### **3.1.LIMITATAION OF SYSTEM**

The main challenge faced in agriculture sector is the lack of knowledge about the changing variations in climate.Each crop has its own suitable climatic features. This can be handled with the help of precise farming techniques.The precision farming not only maintains the productivity of crops but also

increases the yield rate of production.The existing system which recommends crop yield is either hardware-based being costly to maintain, or not easilyaccessible. Despite many solutions that have been recently proposed, there are still open challenges in creating aan application with respect to crop recommendation

### **4. PROPOSED SYSTEM**

Crop production depends on many agricultural parameters.Proposed work is based on the production of crops inprevious years, crops can be recommended to the farmers. This kind of suggestions will make farmer to knowthat whether that particular is yielding a good production in recent years. Production of crops may become lessdue to any crop disease, water problem and many other factors. While considering about the production, farmersmay get knowledge about which crop is in high volume in the market in that year. Based on this farmer can takedecision of trend on crops in recent years. Farmers will be given recommendation by considering the season ofcrop production. Tamilnadu Agriculture Dataset of about 1,20,000 records were taken. It contains fields like cropyear, crop

name, District, Season, Area cultivated and production. Recommendations were given to user based on the production of crops, season when the crops cultivated

#### 4.1 ADVANTAGES OF PROPOSED SYSTEM

The proposed model predicts the crop yield for the data sets of the given region. Integrating agriculture and ML will contribute to more enhancements in the agriculture sector by increasing the yields and optimizing the resources involved. The data from previous years are the key elements in forecasting current performance. The proposed system uses recommender system to suggest the right time for using fertilizers. The methods in the proposed system includes increasing the yield of crops, real-time analysis of crops, selecting efficient parameters, making smarter decisions and getting better yield.

#### 5. MODULE DESCRIPTION

Many harvest yield expectation models have been developed. Bunching methods such as k-means and k-means++ are used to collect data as groupings in order to predict agricultural yield [1]. Tripathy et al. [2]

proposed a methodology for obtaining pesticide executives for crop development via an information mining process. The nature of soil is a fundamental limit for agricultural investigation. In India, there are many different types of soil. Crops are produced based on the soil type in the area. The role of soil in advancing harvest development is discussed [3]. The dirt boundary is investigated using data mining techniques. The JRip, J48, and Naive Bayes techniques are used [4], resulting in more reliable results when dissecting red and black dirt. The impact of agricultural business boundaries on crop executives is investigated in order to improve efficiency [5]. The farming factors are being studied using brain organisations, delicate processing, large data, and fluffy rationale procedures. Pritam Bose [6] developed an SNN model for spatiotemporal analysis and crop evaluation. [7] A programmed framework was constructed to compile data regarding soil nature and weather patterns, using bunching procedures to separate the data and use it by ranchers in crop development. ICT-based communication overcomes any barrier among farmers, such as language barriers. In



this day and age, mobile devices communicate information quickly. Ranchers can use Semantic Web-based Architecture [8] and GIS developments to learn about harvest ideas in a brief amount of time. GIS transmits data about climatic conditions and geographic characteristics. Ranchers can then view this information using any ICT device. GIS and spatial developments can reveal the universe's monetary development [9]. Appropriate processes should be used to extract information from an enormous agriculture data source. Data Mining is an important aspect of the procedures. By using mining, stowed useful information can be retrieved, as well as future forecasts. The information gathered is organized;

### 5.1 Dataset Collection

The dataset comprising the soil specific attributes which are collected for Madurai district tested at soil testing lab, Madurai, Tamil Nadu, India. In addition, similar online sources of general crop data were also used. The crops considered in our model include millet, groundnut, pulses, cotton, vegetables, banana, paddy, sorghum, sugarcane, coriander. Figure 1

gives an analysis of the dataset. The number of instances of each crop available in the training dataset is depicted. The attributes considered where Depth, Texture, Ph, Soil Color, Permeability, Drainage, Water holding and Erosion

ables and 2 area  
yield in different  
work of geospatial

dataset is depicted. The attributes considered where Depth,  
Texture, Ph, Soil Color, Permeability, Drainage, Water  
holding and Erosion.

problem of selecting  
A method to select  
lassifiers has been  
higher accuracy and  
proposed based on  
Using Q statistics,  
ant and accurate  
s which were not  
mble. This measure  
ce and diversity of  
SA (Selection by  
and Diversity) and  
tified. Finally it is  
rs. The paper [8]

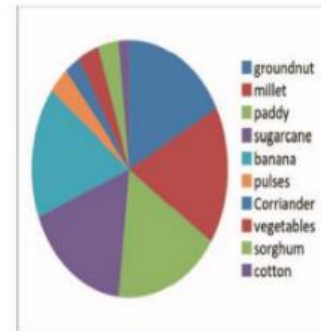


Fig 1 Analysis of dataset with respect to crops(Training data)

The above stated parameters of soil play a major role in the crop's ability to extract water and nutrients from the soil. For crop growth to their fullest potential, the soil must provide a satisfactory environment for it. Soil is the anchor of the roots. The water holding capacity determines the crop's ability to absorb nutrients and other nutrients that are changed into ions, which is the form that the plant can use. Texture determines how porous the soil is and the comfort of air and water movement which is essential to

prevent the plants from becoming waterlogged. Soil texture which affects the soil's ability to hold onto nutrients. The level of acidity or alkalinity (pH) is a master variable which affects the availability of soil nutrients. The activity of microorganisms present in the soil and also the level of exchangeable aluminum can be affected by pH. The water holding and drainage determine the penetration of roots. Hence for the following reasons the above stated parameters are considered for choosing a crop. Ensemble is a data mining model also known as the Committee Methods or Model Combiners, that combine the power of multiple models to acquire greater prediction, efficiency than any of its models could achieve alone. In our system, we use one of the most familiar ensembling technique called Majority Voting technique. In the voting technique any number of base learners can be used. There has to be at least two base learners. The learners are chosen in a way that they are competent to each other yet being complimentary also. Higher the competition higher is the chance of better prediction. But it is necessary for the learners to be complimentary because when one or few

members make an error, the probability of the remaining members correcting this error would be high. Each learner builds itself into a model. The model gets trained using the training data set provided. When a new sample has to be classified, each model predicts the class on its own. Finally, the class which is predicted by majority of the learners is voted to be the class label of the new sample. This method is implemented in Rapid Miner tool (figure 2, 3, 4, 5) depicts the process implemented in Rapid Miner.

	Crop	Year	Season	Crop Area	Production
0	1997	Kharif	Banana	5619	183740.0
1	1997	Kharif	Horse gram	6849	3040.0
2	1997	Kharif	Onion	2813	37188.0
3	1997	Kharif	Sesamum	1598	580.0
4	1997	Kharif	Small millets	63	50.0
..	...	...	...	...	...
537	2013	Whole Year	Sugarcane	1170	121181.0
538	2013	Whole Year	Sweet potato	2	42.0
539	2013	Whole Year	Tapioca	340	10174.0
540	2013	Whole Year	Tobacco	100	159.0
541	2013	Whole Year	Turmeric	1203	6472.0

## 5.2 Preprocess Dataset

It is a technique that is used to convert the raw data set into a clean data set. Prediction model creation We create data into two models: A) Training model B) Testing

The division of the test and train is done in 0.2 and 0.8 that is 20 and 80 percent respectively.

### 5.3 Model evaluation

We apply the machine learning algorithm for testing part and get the accuracy of this model. Prediction This module based on GUI part. we create a web page using bootstrap. The web page like (Nitrogen, Phosphorous, Potassium, PH value, Humidity, Rainfall, Temperature). now we get the data's from user to compare the dataset values. finally it will predict for the Crop and soil to be planted.

## 6. INTERNAL MODULES

### 6.1 Numpy

NumPy is a Python library used for working with arrays. It also has functions for working in domain of linear algebra, fourier transform, and matrices. NumPy was created in 2005 by Travis Oliphant. It is an open source project and you can use it freely. NumPy stands for Numerical Python.

NumPy arrays are stored at one continuous place in memory unlike lists, so processes can access and manipulate them

very efficiently. This behavior is called locality of reference in computer science. This is the main reason why NumPy is faster than lists. Also it is optimized to work with latest CPU architectures.

### 6.2 Pandas

Pandas is a Python library used for working with data sets. It has functions for analyzing, cleaning, exploring, and manipulating data. The name 'Pandas' has a reference to both 'Panel Data', and 'Python Data Analysis' and was created by Wes McKinney in 2008.

Pandas allows us to analyze big data and make conclusions based on statistical theories. Pandas can clean messy data sets, and make them readable and relevant. Relevant data is very important in data science.

### 6.3 Matplotlib

Human minds are more adaptive for the visual representation of data rather than textual data. We can easily understand things when they are visualized. It is better to represent the data through the graph where we can analyze the data more efficiently and make the specific decision according to data

analysis. Before learning the matplotlib, we need to understand data visualization and why data visualization is important.

Graphics provides an excellent approach for exploring the data, which is essential for presenting results. Data visualization is a new term. It expresses the idea that involves more than just representing data in the graphical form (instead of using textual form).

This can be very helpful when discovering and getting to know a dataset and can help with classifying patterns, corrupt data, outliers, and much more. With a little domain knowledge, data visualizations can be used to express and demonstrate key relationships in plots and charts. The static does indeed focus on quantitative description and estimations of data. It provides an important set of tools for gaining a qualitative understanding.

#### 6.4 Keras

Keras is an open-source high-level Neural Network library, which is written in Python is capable enough to run on Theano, TensorFlow, or CNTK. It was developed by

one of the Google engineers, Francois Chollet. It is made user-friendly, extensible, and modular for facilitating faster experimentation with deep neural networks. It not only supports Convolutional Networks and Recurrent Networks individually but also their combination.

It cannot handle low-level computations, so it makes use of the Backend library to resolve it. The backend library act as a high-level API wrapper for the low-level API, which lets it run on TensorFlow, CNTK, or Theano. Initially, it had over 4800 contributors during its launch, which now has gone up to 250,000 developers. It has a 2X growth ever since every year it has grown. Big companies like Microsoft, Google, NVIDIA, and Amazon have actively contributed to the development of Keras. It has an amazing industry interaction, and it is used in the development of popular firms like Netflix, Uber, Google, Expedia, etc.

Focus on user experience has always been a major part of Keras. Large adoption in the industry. It is a multi backend and supports multi-platform, which helps all

the encoders come together for coding. Research community present for Keras works amazingly with the production community. Easy to grasp all concepts. It supports fast prototyping. It seamlessly runs on CPU as well as GPU. It provides the freedom to design any architecture, which then later is utilized as an API for the project. It is really very simple to get started with. Easy production of models actually makes Keras special.

### 6.5 Tensorflow

TensorFlow is a software library or framework, designed by the Google team to implement machine learning and deep learning concepts in the easiest manner. It combines the computational algebra of optimization techniques for easy calculation of many mathematical expressions.

### 6.6 Scikit-learn

Scikit-learn (Sklearn) is the most useful and robust library for machine learning in Python. It provides a selection of efficient tools for machine learning and statistical modeling including classification, regression, clustering and dimensionality reduction via a consistent interface in Python.

## 7. ALGORITHMS USED

### 7.1 Decision tree classifiers

Decision tree classifiers are used successfully in many diverse areas. Their most important feature is the capability of capturing descriptive decision making knowledge from the supplied data. Decision tree can be generated from training sets. The procedure for such generation based on the Set of objects (S), each belonging to one of the classes  $C_1, C_2, \dots, C_k$  is as follows: Step 1. If all the objects in S belong to the same class, for example  $C_i$ , the decision tree for S consists of a leaf labeled with this class. Step 2. Otherwise, let T be some test with possible outcomes  $O_1, O_2, \dots, O_n$ . Each object in S has one outcome for T so the test partitions S into subsets  $S_1, S_2, \dots, S_n$  where each object in  $S_i$  has outcome  $O_i$  for T. T becomes the root of the decision tree and for each outcome  $O_i$  we build a subsidiary decision tree by invoking the same procedure recursively on the set  $S_i$ .

### 7.2 Logistic regression Classifiers

Logistic regression analysis studies the association between a categorical dependent variable and a set of independent

(explanatory) variables. The name logistic regression is used when the dependent variable has only two values, such as 0 and 1 or Yes and No. The name multinomial logistic regression is usually reserved for the case when the dependent variable has three or more unique values, such as Married, Single, Divorced, or Widowed. Although the type of data used for the dependent variable is different from that of multiple regression, the practical use of the procedure is similar.

Logistic regression competes with discriminant analysis as a method for analyzing categorical-response variables. Many statisticians feel that logistic regression is more versatile and better suited for modeling most situations than is discriminant analysis. This is because logistic regression does not assume that the independent variables are normally distributed, as discriminant analysis does.

This program computes binary logistic regression and multinomial logistic regression on both numeric and categorical independent variables. It reports on the regression equation as well as the goodness of fit, odds ratios, confidence limits,

likelihood, and deviance. It performs a comprehensive residual analysis including diagnostic residual reports and plots. It can perform an independent variable subset selection search, looking for the best regression model with the fewest independent variables. It provides confidence intervals on predicted values and provides ROC curves to help determine the best cutoff point for classification. It allows you to validate your results by automatically classifying rows that are not used during the analysis.

### 7.3 Naïve Bayes

The naive bayes approach is a supervised learning method which is based on a simplistic hypothesis: it assumes that the presence (or absence) of a particular feature of a class is unrelated to the presence (or absence) of any other feature.

Yet, despite this, it appears robust and efficient. Its performance is comparable to other supervised learning techniques. Various reasons have been advanced in the literature. In this tutorial, we highlight an explanation based on the representation bias. The naive bayes classifier is a linear

classifier, as well as linear discriminant analysis, logistic regression or linear SVM (support vector machine). The difference lies on the method of estimating the parameters of the classifier (the learning bias). While the Naive Bayes classifier is widely used in the research world, it is not widespread among practitioners which want to obtain usable results. On the one hand, the researchers found especially it is very easy to program and implement it, its parameters are easy to estimate, learning is very fast even on very large databases, its accuracy is reasonably good in comparison to the other approaches. On the other hand, the final users do not obtain a model easy to interpret and deploy, they do not understand the interest of such a technique.

Thus, we introduce in a new presentation of the results of the learning process. The classifier is easier to understand, and its deployment is also made easier. In the first part of this tutorial, we present some theoretical aspects of the naive bayes classifier. Then, we implement the approach on a dataset with Tanagra. We compare the obtained results (the parameters of the model) to those obtained with other linear

approaches such as the logistic regression, the linear discriminant analysis and the linear SVM. We note that the results are highly consistent. This largely explains the good performance of the method in comparison to others. In the second part, we use various tools on the same dataset (Weka 3.6.0, R 2.9.2, Knime 2.1.1, Orange 2.0b and RapidMiner 4.6.0). We try above all to understand the obtained results.

#### 7.4 Random Forest

Random forests or random decision forests are an ensemble learning method for classification, regression and other tasks that operates by constructing a multitude of decision trees at training time. For classification tasks, the output of the random forest is the class selected by most trees. For regression tasks, the mean or average prediction of the individual trees is returned. Random decision forests correct for decision trees' habit of overfitting to their training set. Random forests generally outperform decision trees, but their accuracy is lower than gradient boosted trees. However, data characteristics can affect their performance. The first algorithm for random decision forests was created in 1995 by Tin

KamHo[1] using the randomsubspace method, which, in Ho&'ss formulation, is a way to implement the 'sstochastic discrimination'sapproach to classification proposed by Eugene Kleinberg. An extension of the algorithm was developed by Leo Breiman and Adele Cutler, who registered'sRandomForests's as a trademark in 2006 (as of 2019, owned by Minitab, Inc.). The extension combines Breiman&'ss'bagging's idea and random selection of features, introduced first by Ho[1] and later independently by Amit and Geman[13] in order to construct a collection of decision trees with controlled variance. Random forests are frequently used as 'sblackbox's models in businesses, as they generate reasonable predictions across a wide range of data while requiring little configuration.

## 7.5 SVM

In classification tasks a discriminant machine learning technique aims at finding, based on an independent and identically distributed (iid) training dataset, a discriminant function that can correctly predict labels for newly acquired instances. Unlike generative machine

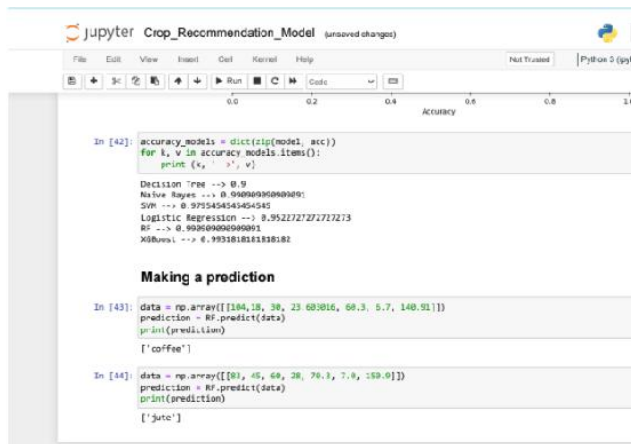
learning approaches, which require computations of conditional probability distributions, a discriminant classification function takes a data point  $x$  and assigns it to one of the different classes that are apart of the classification task. Less powerful than generative approaches, which are mostly used when prediction involves outlier detection, discriminant approaches require fewer computational resources and less training data, especially for a multidimensional feature space and when only posterior probabilities are needed. From a geometric perspective, learning a classifier is equivalent to finding the equation for a multidimensional surface that best separates the different classes in the feature space.

SVM is a discriminant technique, and, because it solves the convex optimization problem analytically, it always returns the same optimal hyperplane parameter—in contrast to genetic algorithms (GAs) or perceptrons, both of which are widely used for classification in machine learning. For perceptrons, solutions are highly dependent on the initialization and termination criteria. For a specific kernel that transforms the data from the input space



to the feature space, training returns uniquely defined SVM model parameters for a given training set, whereas the perceptron and GA classifier models are different each time training is initialized. The aim of Gasand perceptrons is only to minimize error during training, which will translate into several hyperplanes' meeting this requirement.

## 8. OUTPUT RESULTS



```

Jupyter Crop_Recommendation_Model (unsaved changes)
Python 3 (ipy)

In [42]: accuracy_models = dict(zip(model_acc))
for k, v in accuracy_models.items():
    print(k, ' -> ', v)

Decision Tree -> 0.9
Naive Bayes -> 0.9909090909090909
SVM -> 0.9754545454545454
Logistic Regression -> 0.9522727272727273
RF -> 0.9909090909090909
XGBowl -> 0.9931818181818182

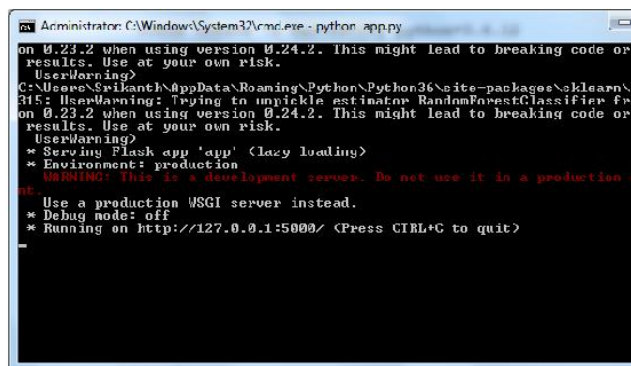
Making a prediction

In [43]: data = np.array([[104, 18, 30, 23, 683016, 68.3, 5.7, 140.9]])
prediction = RF.predict(data)
print(prediction)

['coffee']

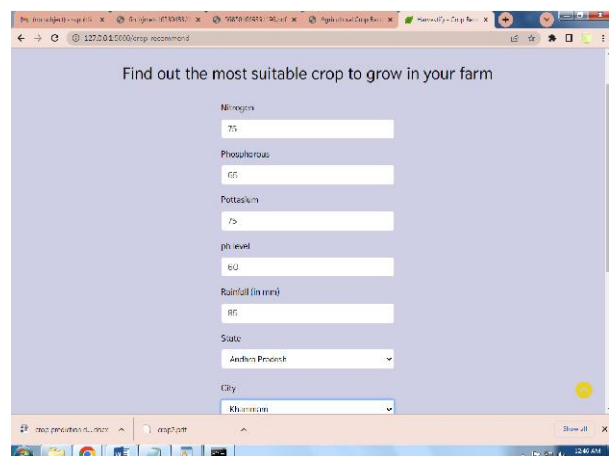
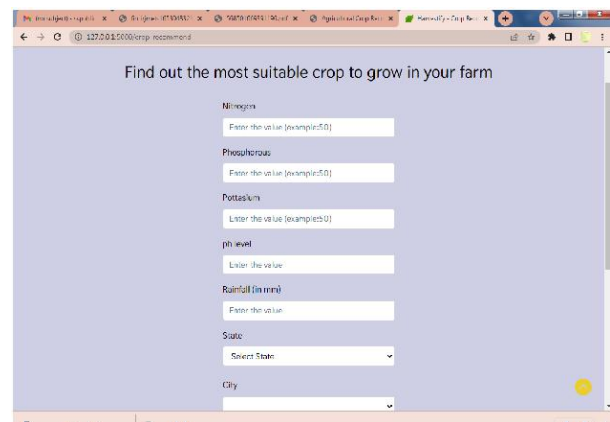
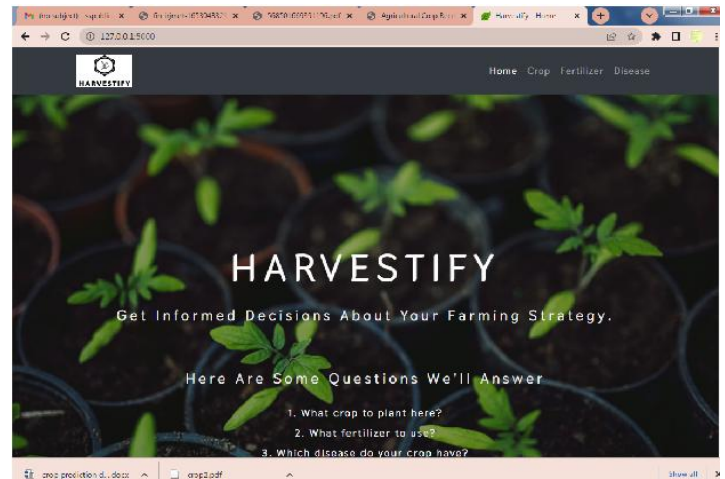
In [44]: data = np.array([[81, 45, 64, 38, 76.3, 7.8, 159.0]])
prediction = RF.predict(data)
print(prediction)

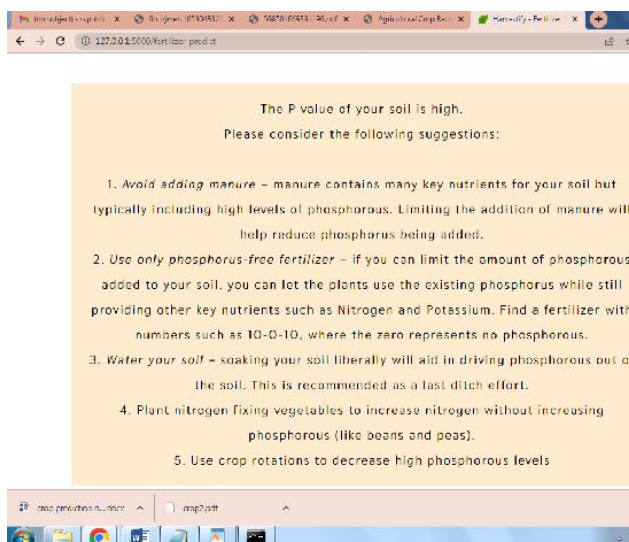
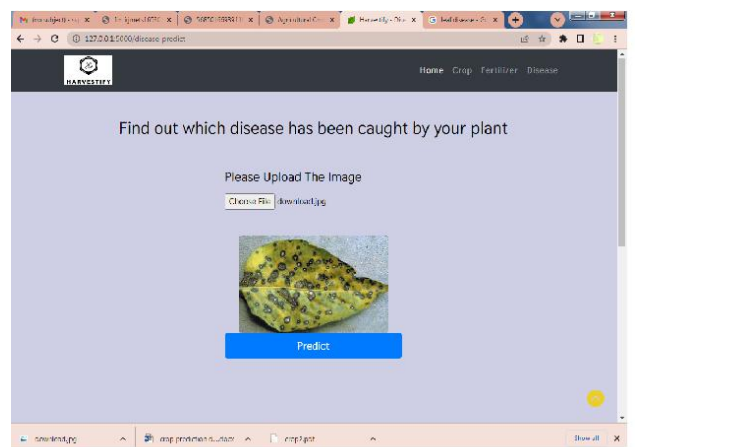
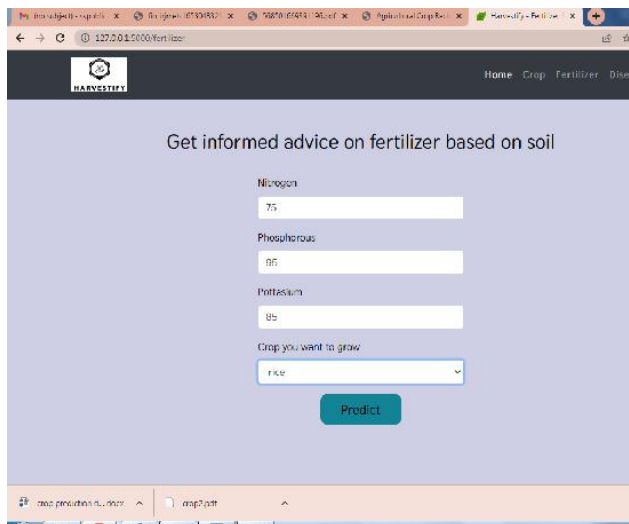
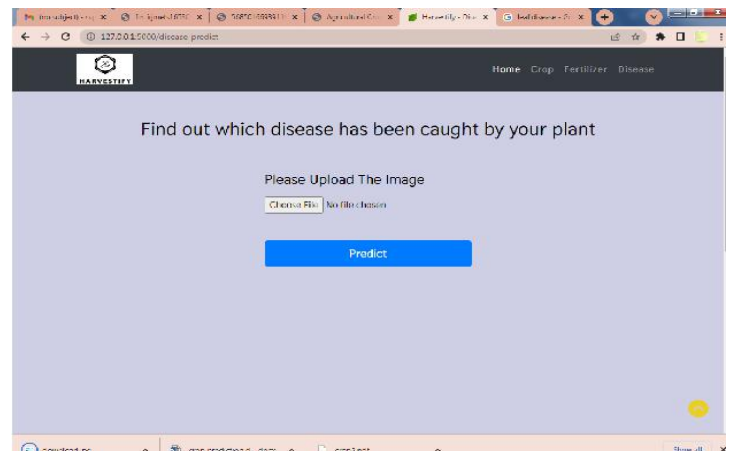
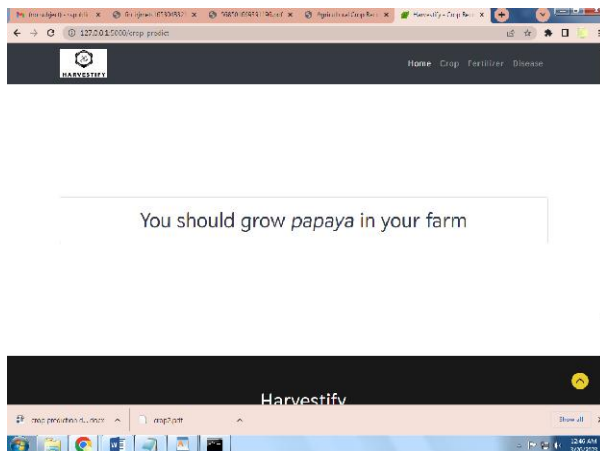
['jute']
    
```



```

Administrator: C:\Windows\System32\cmd.exe - python app.py
on 0.23.2 when using version 0.24.2. This might lead to breaking code or
results. Use at your own risk.
UserWarning:
C:\Users\Selkenth\AppData\Roaming\Python\Python36\site-packages\sklearn\
utils\validation.py:115: UserWarning: Trying to unpickle estimator RandomForestClassifier fr
on 0.23.2 when using version 0.24.2. This might lead to breaking code or
results. Use at your own risk.
UserWarning:
* Serving Flask app "app" (lazy loading)
* Environment: production
WARNING: This is a development server. Do not use it in a production
setting.
Use a production WSGI server instead.
* Debug mode: off
* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
    
```





## 9. CONCLUSION

In this paper, significance of management of crops was studied vastly. Farmers need assistance with recent technology to grow their crops. Proper prediction of crops can be informed to agriculturists in time basis. Many Machine Learning techniques have been used to analyze the agriculture parameters. Some of the techniques indifferent aspects of agriculture are studied

by a literature study. Blooming Neural networks, Soft computing techniques plays significant part in providing recommendations. Considering the parameter like production and season, more personalized and relevant recommendations can be given to farmers which makes them to yield good volume of production.

## 10. REFERENCES

- [1] Shreya S. Bhanose, Kalyani A. Bogawar (2016) “Crop And Yield Prediction Model”, International Journal of Advance Scientific Research and Engineering Trends, Volume 1, Issue 1, April 2016
- [2] Tripathy, A. K., et al. (2011) ‘sData mining and wireless sensor network for agriculture pest/disease predictions.’s Information and Communication Technologies (WICT), 2011 World Congress on. IEEE.
- [3] Ramesh Babu Palepu (2017) ” An Analysis of Agricultural Soils by using Data Mining Techniques”, International Journal of Engineering Science and Computing, Volume 7 Issue No. 10 October.
- [4] Rajeswari and K. Arunesh (2016) “Analysing Soil Data using Data Mining Classification Techniques”, Indian Journal of Science and Technology, Volume 9, May.
- [5] A.Swarupa Rani (2017), “The Impact of Data Analytics in Crop Management based on Weather Conditions”, International Journal of Engineering Technology Science and Research, Volume 4, Issue 5, May.
- [6] Pritam Bose, Nikola K. Kasabov (2016), “Spiking Neural Networks for Crop Yield Estimation Based on Spatiotemporal Analysis of Image Time Series”, IEEE Transactions On Geoscience And Remote Sensing.
- [7] Priyanka P.Chandak (2017),” Smart Farming System Using Data Mining”, International Journal of Applied Engineering Research, Volume 12, Number 11.
- [8] Vikas Kumar, Vishal Dave (2013), “KrishiMantra: Agricultural Recommendation System”, Proceedings of the 3rd ACM Symposium on Computing for Development, January.
- [9] SavaeLatu (2009), ”Sustainable Development : The Role Of Gis And Visualisation”, The Electronic Journal on Information Systems in Developing Countries, EJISDC 38, 5, 1-17.

[10] NasrinFathima.G (2014), “Agriculture Crop Pattern Using Data Mining Techniques”, International Journal ofAdvanced Research in Computer Science and Software Engineering, Volume 4, May.

[11] Ramesh A.Medar (2014), ”A Survey on Data Mining Techniques for Crop Yield Prediction”, InternationalJournal of Advance Research in Computer Science and Management Studies, Volume 2, Issue 9, September.

[12] ShakilAhamed.A.T.M, NavidTanzeem Mahmood (2015),” Applying data mining techniques to predict annualyield of major crops and recommend planting different crops in different districts in Bangladesh”, ACIS 16<sup>th</sup>International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/DistributedComputing (SNPD),IEEE,June.

[13] Shreya S.Bhanose (2016),”Crop and Yield Prediction Model”, International Journal of Advence ScientificResearch and Engineering Trends, Volume 1,Issue 1,ISSN(online) 2456- 0774,April.

[ 14] AgajiIorshase, OnyekiIdokoCharles,”A Well-Built Hybrid Recommender System for Agricultural Products inBenue State of Nigeria”, Journal of Software Engineering and Applications,2015,8,581-589

[15] G. Adomavicius and A. Tuzhilin(2005), “Toward the Next Generation of Recommender Systems: A Survey ofthe State-of-theArt and Possible Extensions,” IEEE Trans. Knowledge and Data Eng., vol. 17, no. 6, pp. 734-749,June.

[16] Avinash Jain, Kiran Kumar (2016),”Application of Recommendation Engines in Agriculture”, InternationalJournal of Recent Trends in Engineering & Research, ISSN: 2455-1457.

[ 17] Kiran Shinde (2015),”Web Based Recommendation System for farmers”, International Journal on Recent andInnovation Trends in Computing and Communication, Volume 3,Issue 3, ISSN:2321- 8169,March.

[18] Konstantinos G. Liakos, “ Machine Learning in Agriculture: A Review”,

Sensors 2018, 18,  
2674;doi:10.3390/s18082674

[19] S.Vaishnavi, M.Shobana, N Geethanjali,  
Dr.S.Karthik, “Data Mining: Solving the  
Thirst of Recommendations toUsers”, IOSR  
Journal of Computer Engineering (IOSR-  
JCE), Vol.16, no.6, 2014