

# **CROP YIELD PREDICTION USING DEEP LEARNING**

## **A PROJECT REPORT**

*Submitted by*

<b>CB.EN.U4ECE18107</b>	<b>ARUNESWARI S</b>
<b>CB.EN.U4ECE18110</b>	<b>BOMMINENI RAMAKRISHNA</b>
<b>CB.EN.U4ECE18112</b>	<b>CHANDHANA S</b>
<b>CB.EN.U4ECE18136</b>	<b>NIRALI M DAVE</b>
<b>CB.EN.U4ECE18144</b>	<b>PRATHIKSHA R R</b>

*Under the guidance of*

**Dr. ARAVINTH J**

*in partial fulfillment of the requirements for the award of the  
degree of*



**BACHELOR OF TECHNOLOGY**

**IN**

**ELECTRONICS AND COMMUNICATION ENGINEERING**

**AMRITA SCHOOL OF ENGINEERING**

**AMRITA VISHWA VIDYAPEETHAM**

**COIMBATORE 641112**

**May 2022**

# AMRITA VISHWA VIDYAPEETHAM

AMRITA SCHOOL OF ENGINEERING, COIMBATORE, 641112



## BONAFIDE CERTIFICATE

This is to certify that the project report entitled  
“CROP YIELD PREDICTION USING DEEP LEARNING”.

submitted by

<b>CB.EN.U4ECE18107</b>	<b>ARUNESWARI S</b>
<b>CB.EN.U4ECE18110</b>	<b>BOMMINENI RAMAKRISHNA</b>
<b>CB.EN.U4ECE18112</b>	<b>CHANDHANA S</b>
<b>CB.EN.U4ECE18136</b>	<b>NIRALI M DAVE</b>
<b>CB.EN.U4ECE18144</b>	<b>PRATHIKSHA R R</b>

in partial fulfillment of the requirements for the award of the **Degree of Bachelor of Technology** in **ELECTRONICS AND COMMUNICATION ENGINEERING** is a Bonafide record of the work carried out under my guidance and supervision at Amrita School of Engineering, Coimbatore.

Project Advisor  
Name: **Dr. Aravinth J**  
Designation: **Assistant Professor**  
(Selection Grade)

Project Coordinator  
Name: **M. E. Harikumar**  
Designation: **Assistant Professor**

Chairperson ECE  
Dr. M. Jayakumar

The project was evaluated by us on:

Internal Examiner

External Examiner

## **ACKNOWLEDGEMENT**

We would like to thank **Dr. Aravinth J**, Assistant Professor, Department of Electronics and Communication Engineering, for his invaluable direction, patience, and assistance during the entire process, which allowed us to complete this project successfully and on schedule. We appreciate our panel members' intelligent remarks, challenging questions, and support, which helped us improve our project. We'd also like to express our gratitude to all of the teaching and non-teaching staff who have mentored and supported us. We also like to thank our parents and friends for their continual encouragement and support.

## **ABSTRACT**

The use of remote sensing in smart farming is gaining popularity around the world. Due to the rising global demand for food grains, assessing yield before actual production is critical in developing policies and making decisions in the agricultural production system. In this project, the aim is to develop a deep learning model which will predict an individual farm's crop yield using remotely sensed satellite images. The study focuses on crop yield estimation of wheat, corn and soybean in farms in the chosen study area in KBS-LTER, Michigan, United States. The satellite data used is collected from Landsat 5 and 7, and the following 7 spectral bands were retrieved: blue, red, green, near-infrared, short wave infrared 1, short wave infrared 2, and thermal infrared. The seven spectral bands were investigated for the pooled data of 11 years (2001 to 2012, excluding 2002) years. Four widely used vegetation spectral indices were calculated from the spectral bands for the 11 years. These indices were (a) Normalized Difference Vegetation Index (NDVI), (b) Soil Adjusted Vegetation index (SAVI), (c) Green Red NDVI (GRNDVI), and (d) Two band Enhanced Vegetation Index (EVI2). The dataset used to train the model is developed from the original KBS LTER geo-referenced annual crop yields data points. To improve the accuracy, historical images (Landsat 5 and 7) were collected from GEE based on the plantation and harvest dates of the crop. Each farm under study is cropped out into a cluster of adjacent pixels to increase the dataset. Further augmentation methods were implemented to increase the dataset, which eradicated the problem of overfitting of the model when it was trained with 252 images only. These methodologies increased the dataset to 23,072. The yield corresponding to each 4-pixel image is annotated.

Two Artificial Neural Network models were developed and trained using the 7 spectral dataset and 4 vegetation indices dataset that was created, which predicts the farm level crop yield. The accuracy of the Pool-Band model was 97.34% compared to the Vegetation index model having an accuracy of 88.8%. To improve the accuracy of the

model, feature engineering is performed. The methods used are: 1. Log transform of yield values 2. Adding historical data.

Further year-wise (11 models) and crop-wise analysis (3 models) were performed with respect to 7-spectral bands as input features.

## TABLE OF CONTENTS

CHAPTER NO.	TITLE	PAGE NO.
	<b>ABSTRACT</b>	i
	<b>LIST OF ABBREVIATIONS</b>	iii
	<b>LIST OF FIGURES</b>	v
	<b>LIST OF TABLES</b>	viii
<b>1</b>	<b>INTRODUCTION</b>	1
	1.1 INTRODUCTION	2
<b>2</b>	<b>LITERATURE SURVEY</b>	5
	2.1 LITERATURE SURVEY	6
	2.1.1 REGRESSION BASED MODEL	6
	2.1.2 SIMULATION BASED MODEL	7
	2.1.3 DEEP LEARNING BASED MODEL	8
<b>3</b>	<b>PROPOSED WORK</b>	10
	3.1 WORKFLOW	11
	3.2 STUDY AREA	12
	3.2.1 CROPS IN STUDY AREA	14
	3.3 RETRIEVAL OF SATELLITE DATA	15
	3.3.1 HISTORICAL DATA RETRIEVAL	19
	3.4 COMPUTATION OF CROP YIELD	20
	3.4.1 LOG TRANSFORMATION TECHNIQUE	20
	3.5 EXTRACTING BAND VALUES	21
	3.6 EXTRACTING VEGETATION INDICES	22
	3.7 METHODS TO INCREASE DATASET	24
	3.7.1 AUGMENTATION	25
	3.8 ANNOTATION	25
	3.9 DEEP LEARNING	26

	3.10 NEURAL NETWORK	26
	3.11 ANN MODEL - SPECTRAL BANDS	29
	3.12 ANN MODEL - VEGETATION INDICES	31
<b>4</b>	<b>RESULTS AND DISCUSSIONS</b>	33
	4.1 SOFTWARE PLATFORMS	34
	4.2 RESULTS	34
	4.2.1 DATASET PREPARATION	34
	4.2.2 PERFORMANCE METRICS	36
	4.2.3 ANALYSIS OF POOL SPECTRAL BANDS	37
	4.2.4 ANALYSIS OF POOL VEGETATION INDICES	38
	4.2.5 CROP-WISE ANALYSIS	40
	4.2.6 YEAR-WISE ANALYSIS	41
	4.3 DISCUSSIONS	47
<b>5</b>	<b>CONCLUSIONS AND FUTURE WORK</b>	49
	<b>REFERENCES</b>	51
	<b>PUBLICATIONS</b>	54

## LIST OF ABBREVIATIONS

ABBREVIATION	EXPANSION	PAGE NO
SVM	Support Vector Machine	3
RF	Random forest	3
ERT	Extremely Randomized Trees	3
NDVI	Normalized difference vegetation index	4
SAVI	Soil-Adjusted Vegetation Index	4
GRNDVI	Green Red Normalized difference vegetation index	4
EVI	Enhanced vegetation index	6
GEE	Google earth engine	17
KBS LTER	Kellogg Biological Station - Long-term ecological research	3
GPS	Global Positioning System	2
ETM+	Enhanced Thematic Mapper Plus	3
TM	Thematic Mapper	3
NIR	Near infrared	3
ANN	Artificial Neural Network	4
LAI	Leaf Area Index	6
DL	Deep Learning	6
MODIS	Moderate Resolution Imaging Spectroradiometer	6
ML	Machine Learning	6
OLI	Operational Land Imager	7
PAR	Photosynthetically active radiation	7
SAFY	Simple Algorithm For Yield	7
SLA	Specific leaf area	8



ELUE	Effective light use efficiency	8
CNN	Convolutional Neural Network	8
LSTM	Long Short-Term Memory	8
HLS	Harmonized Landsat Sentinel	8
LASSO	Least absolute shrinkage and selection operator	9
NASA	National Aeronautics and Space Administration	8
RMSE	Root-mean-square error	9
BPNN	Back-propagation neural network	9
PVI	Perpendicular vegetation index	9
KML	Keyhole Markup Language	13
MSS	Multispectral Scanner	16
GIS	Geographic Information System	17
ReLU	Rectified Linear Unit	29
MAPE	Mean Absolute Percentage Error	31

## LIST OF FIGURES

FIGURE NO	TITLE	PAGE NO
3.1	Flowchart of proposed model	11
3.2	Geographic data plotted in Google Earth Pro for the year 2001	14
3.3	Cloud cover of 93% (Landsat 5)	18
3.4	Stripes in image (Landsat 7)	18
3.5	Farms Visualization in RGB Format	18
3.6	After plantation for the year 2001 (Wheat)	19
3.7	Before harvest for the year 2001 (Wheat)	19
3.8	Augmentation	25
3.9	Neural Network	26
3.10	Computation in neuron	28
4.1	Original yield and Predicted yield - All bands	37
4.2	Model Loss - All bands	37
4.3	Original yield vs Predicted yield - All bands	37
4.4	Original yield and Predicted yield -Vegetation Indices	38
4.5	Model Loss - Vegetation Indices	38
4.6	Original yield vs Predicted yield-Vegetation Indices	39

4.7	Crop yield and Predicted yield for Corn	40
4.8	Crop yield and Predicted yield for Soybean	40
4.9	Crop yield and Predicted yield for Wheat	41
4.10	Original yield and Predicted yield for the year 2001	41
4.11	Original yield vs Predicted yield for the year 2001	41
4.12	Original yield and Predicted yield for the year 2003	42
4.13	Original yield vs Predicted yield for the year 2003	42
4.14	Original yield and Predicted yield for the year 2004	42
4.15	Original yield vs Predicted yield for the year 2004	42
4.16	Original yield and Predicted yield for the year 2005	43
4.17	Original yield vs Predicted yield for the year 2005	43
4.18	Original yield and Predicted yield for the year 2006	43
4.19	Original yield vs Predicted yield for the year 2006	43
4.20	Original yield and Predicted yield for the year 2007	44
4.21	Original yield vs Predicted yield for the year 2007	44
4.22	Original yield and Predicted yield for the year 2008	44
4.23	Original yield vs Predicted yield for the year 2008	44
4.24	Original yield and Predicted yield for the year 2009	45
4.25	Original yield vs Predicted yield for the year 2009	45
4.26	Original yield and Predicted yield for the year 2010	45

4.27	Original yield vs Predicted yield for the year 2010	45
4.28	Original yield and Predicted yield for the year 2011	46
4.29	Original yield vs Predicted yield for the year 2011	46
4.30	Original yield and Predicted yield for the year 2012	46
4.31	Original yield vs Predicted yield for the year 2012	46

## LIST OF TABLES

<b>TABLE NO</b>	<b>TITLE</b>	<b>PAGE NO</b>
3.1	Dataset description	13
3.2	Planting and Harvesting period of the crops	15
3.3	Year-wise Harvest Dates of the crop	15
3.4	Landsat image retrieval dates	20
3.5	Multispectral Bands	21
3.6	Hyperparameter tuning – Spectral Band model	31
3.7	Hyperparameter tuning – Vegetation Indices model	32
4.1	Size of dataset before augmentation	35
4.2	Size of dataset after augmentation	35
4.3	Feature size	35
4.4	Performance of pool spectral band model	37
4.5	Performance of pool vegetation indices model	39
4.6	Crop-wise analysis	41
4.7	Year-wise analysis	46

# **CHAPTER 1**

## **INTRODUCTION**

## 1.1 INTRODUCTION

With the ever-increasing global population, adequate crop production is essential. Therefore, agricultural systems constitute a pivotal economic sector worldwide. Sustainable rural systems are very much dependent upon the healthy development of agriculture. Remote sensing has proved to be game-changing for the agricultural sector, as it is one of the backbones of precision in agriculture. In remote sensing, multispectral and hyperspectral satellite images play a major role in crop management and estimation of crop yields, their ability to represent crop growth conditions on the spatial and temporal scale is remarkable. These images can describe crop development using a variety of vegetation indices. Many approaches have been developed to translate remote sensing data into crop yields, and several reviews of such approaches exist. Farmers have been able to convert from conventional farming methods to data-driven decision-making because of an increase in data-producing equipment and sensors, which has been a consistent trend in agriculture. This is referred to as "smart farming." [1][2]

Some of these approaches have faced several problems in estimating crop yield. One of these problems is the scarcity of remote sensing data suitable for use in crop management because of climatic conditions such as clouds and the scarcity of crop yield patterns for the individual farms. Most of the time climatic conditions and low temporal resolution are the main obstacles that prevent decision-makers from using remote sensing data to map crops and estimate crop yields. Previously, Statistical methods that use regression approaches had been adopted to predict the crop yield based on numerous climatic factors such as rainfall, temperature, drought, etc. [3] Regression-based models presume that crop yield can be calculated using only a few independent factors. These models may lead to erroneous results due to the assumptions made.[4]. Crop yields can be estimated county-wise, state-wise, and farm-wise for a large area of study. Hence, Remote sensing and global positioning systems (GPS) methodologies can be incorporated for better results [5].

The initial methods of using remote sensing to estimate the crop yield was achieved using regression equations.[6]. Traditional machine learning algorithms have been used along

with remote sensing data to carry out various tasks such as crop classification, crop yield prediction, and weed detection. Feature extraction is the initial step in implementing any machine learning algorithm. The process of optimal feature extraction from data is difficult using traditional methods, hence the development and training of multilayer algorithms (Deep Learning) will produce better results. It was noted in [7], that the corn yield estimation is more accurate using Deep Learning in contrast to machine learning algorithms such as SVM, RF, and ERT.

This project aims to develop a deep learning model (Artificial Neural Network) to estimate the crop yield of a farm. Most of these studies have focused on regional/state/county level analysis [7],[8]. It is more beneficial for farmers and agricultural researchers to make quality decisions regarding crop management if the yield estimation is done at the field level rather than the county level. Hence, a detailed study of crop yield prediction at the farm level has been implemented in this work. Our study focuses on predicting the farm-level crop yield of soybean, corn, and wheat in farms under the chosen study area in Michigan, United States using Landsat 5 and Landsat 7 satellite images which are atmospherically corrected surface reflectance from the Landsat 7 ETM+ sensor or the Landsat 5 TM sensor. Seven spectral bands were extracted for the analysis namely: Band 1 (blue) surface reflectance, Band 2(green) surface reflectance, Band 3 (red) surface reflectance, Band 4 (near infrared) surface reflectance, Band 5 (shortwave infrared 1) surface reflectance, Band 6 Thermal Infrared, Band 7 (shortwave infrared 2) surface reflectance and crop yield data provided from KBS LTER for the years 2001-2012 (excluding 2002). These contribute to the input features which were fed to the model.

Vegetation indices were used to draw a comparison with spectral bands, on what can be used as a better input feature for crop yield prediction model. The spectral bands such as red, green, blue band and near-infrared band (NIR) band have been used to monitor vegetation, crop growth, crop health, soil moisture and crop yield[19]. The vegetation indices are a result of various associations between the vegetation such as a crop land and the electromagnetic spectral bands such as red and near infrared bands[20].



The vegetation indices used in the study are NDVI, SAVI, EVI-2, GRNDVI, which will be discussed in section 3.6.

Visualizing the yield data on Google earth pro led to the conclusion that there are 24 farms under study in the area of Michigan. A detailed workflow of the dataset creation and collection is explained in section 3.1. Further pre-processing of the satellite images is done prior to annotation which will be explained in chapter 3.

For the purpose of estimating yield of a farm, two deep learning regression ANN models were developed. These ANN models were trained, one with input features being the 7-spectral bands, and the other with input features being the vegetation indices. The results of these models are compared and analyzed.

## **CHAPTER 2**

### **LITERATURE SURVEY**

## **2.1 LITERATURE SURVEY**

Over the years, research has revealed various methods of yield estimation techniques which are:

- Empirical-based (Regression-based) models
- Process-based (Simulation-based) models
- Models based on deep neural networks

### **2.1.1 REGRESSION BASED MODEL**

MODIS satellite image data are being used to find the corn yield estimation for the state of Iowa in the United States by N Kim et.al [7]. They have summarized the dataset with 3 major categories of parameters, the first being the remote sensing parameters like NDVI, EVI, and LAI, the second being the climatic parameters and the third being the yield of the corn. This dataset is created for each of the counties of the Iowa state for the cropland pixels. The machine learning models that have been used are SVM (support vector machine), random forest, and extremely randomized trees. Overall, it was concluded that the DL method showed the highest accuracy. It is observed that the Machine learning techniques such as SVM, RF, and ERT have an overfitting problem, which occurs when a model is very complex with many parameters and shows a poor predictive performance by overreacting to minor fluctuations in the dataset. This indicates that deep learning models are much more suitable than the SVM, RF, and ERT models when it comes to complex datasets of crop yield prediction.

For instance, in Arumugam et al.'s [9] work, machine learning technique (Gradient Boosted Regression) was implemented for estimating rice yields for the Kharif season from 2003 to 2017 at 500 m spatial resolution in India. The crop mask for rice was taken from the Moderate Resolution Imaging Spectroradiometer (MODIS) multispectral rice classification. This was used to classify one crop over the other. LAI (leaf area index), is a parameter used for the ML regression model. Using the gradient boosted regression model, the yield was predicted.

Gaussian process regression, support vector regression, boosted regression trees, and Random Forest regression were implemented to predict the yield values by Aghighi et al. [10]. Dataset was obtained from Landsat-OLI and NDVI (Vegetation Index) was retrieved. The performance of boosted regression trees technique was better than the others, whose correlation coefficient value was greater than 0.87 for all years. Then followed by random forest regression and Gaussian process regression. Support vector regression did not show a good performance.

Adeniyi et al. [18], implemented Wheat yield prediction by utilizing the time series vegetation indices from Landsat 8 image data for the years 2013 to 2019 in Jász-Nagykun-Szolnok County present in central Hungary. The VI used were Normalized Difference Vegetation Index (NDVI) and Soil Adjusted Vegetation Index (SAVI) and among these two a comparative study was evaluated to find the most accurate VI that predicts the wheat prediction. This was implemented by using a linear regression model that draws a relationship between the wheat yield and the two Vegetation indices. Results from Validation concluded that the SAVI forecast model performed more accurately in comparison to the NDVI forecast model.

Regression-based models presume that crop yield can be predicted by considering only variables that are all independent of one another. This had various shortcomings which led to the usage of simulation-based models.

### **2.1.2 SIMULATION-BASED MODEL**

Simulation-based models take into account the mechanisms of photosynthesis, respiration, crop phenology, the daily photosynthetically active radiation absorbed by the canopy (PAR), etc which are dependent on each other and hence represent a real crop.

A maize yield prediction is implemented on a regional scale by Marjorie Battude et.al[17]. The simulation model used here is the Simple Algorithm For Yield estimates (SAFY) model trained using remote sensing data from Formosat-2, SPOT4-Take5, Landsat-8, and Deimos-1. To increase the accuracy, a modified SAFY model was

developed that also considers the seasonal variation of specific leaf area (SLA) and effective light use efficiency (ELUE).

Since the simulation-based models require a high number of agro-environmental input factors, they can only be used at a local scale. Hence, their application is limited. Therefore, there was a shift towards deep learning-based crop yield prediction models.

### **2.1.3 DEEP LEARNING MODEL**

The process of optimal feature extraction from data is difficult using traditional methods, hence the development and training of multilayer algorithms were introduced. Deep learning algorithms help in better feature extraction properties which help in better crop yield prediction results.

Kuwata et al. [11], implemented a deep learning model to estimate corn yield in Illinois to extract features that the crop growth depends on. Three Algorithms were implemented namely SVR (Support Vector Machine), a single InnerProductLayer, and two InnerProductLayer neural networks, and the highest accuracy was achieved by a deep learning model with two InnerProductLayer.

A county-level analysis and field level analysis has been implemented by Ghazaryan et al. [12], to predict the crop yield of maize and soybean in the United States, using MODIS-based surface reflectance, land surface temperature, and evapotranspiration time series as inputs for county-level analysis and NASA's Harmonized Landsat Sentinel-2 (HLS) product as inputs for field-level analysis. A deep learning model was created using 3D CNN and CNN followed by LSTM.

For instance, in Engen et al. [13]'s work, crop yield prediction is being implemented using Sentinel-2 satellite images, weather data, farm data, grain delivery data, and cadastre-specific data by developing a deep hybrid neural network model consisting of convolutional layers and recurrent neural networks so as to train the multi-temporal data.

A 3D CNN model was developed for soybean yield prediction with the use of satellite imagery by Russello et al. [14] It has been concluded that the 3D CNN model that was developed outperform the existing machine learning methods due to the leverage of spatial, spectral, and temporal dimension of the remote sensing image.

In[15], A long short-term memory (LSTM) model is developed that integrates heterogeneous crop phenology, meteorology, and remote sensing data to evaluate an estimate of county-level corn yield in the US Corn Belt. The LSTM model surpassed the least absolute shrinkage and selection operator (LASSO) regression and random forest (RF) models for estimating county-level corn yield due to its capability to learn patterns from spatial, spectral, and temporal input features.

The study of soybean yield prediction is done in Lauderdale County, Alabama, USA by Terliksiz et al. [8]. A 3D CNN model is being implemented. The satellite data is retrieved from NASA's MODIS land products surface reflectance, and land surface temperature. The evaluation metric adopted is RMSE so as to compare it with other methods.

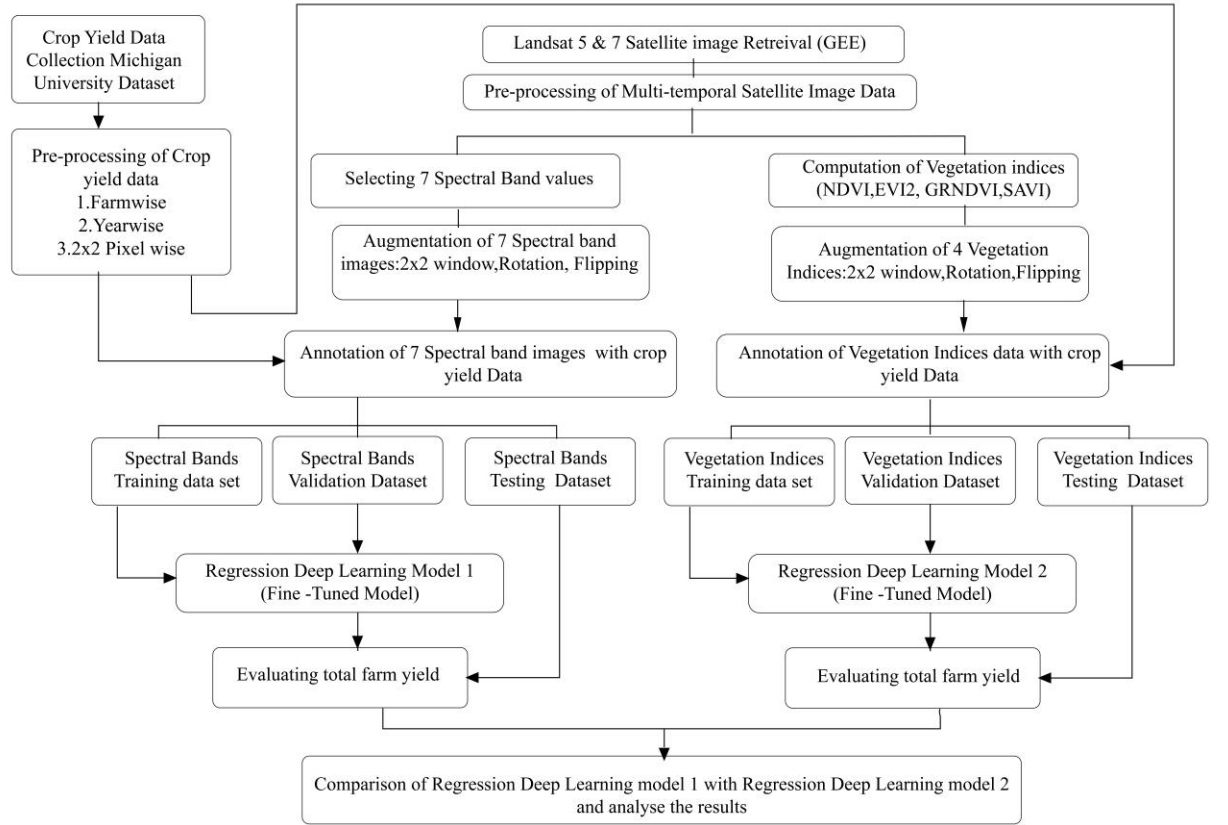
For instance, in Panda et al. [16]'s work, A corn yield prediction model was developed using Back-propagation Neural Network (BPNN) by using indices such as normalized difference vegetation index (NDVI), green vegetation index (GVI), soil adjusted vegetation index (SAVI), perpendicular vegetation index (PVI) for a duration of 3 years. The BPNN model that was developed included 16 models (4 indices \* 4 years including the data from the pooled years). The most accurate BPNN model was PVI grid images by using mean and standard deviation.

## **CHAPTER 3**

### **PROPOSED WORK**

### 3.1 WORKFLOW

The proposed model is illustrated as shown in Fig 3.1. This deals with prediction of crop yield of the farms under study for the years 2003 to 2012 and 2001 from satellite images retrieved from Landsat 5 and Landsat 7, with the help of a regression based deep learning model.



**Fig 3.1** Flowchart of proposed model

As shown in the flowchart Fig 3.1, crop yield ground data of farms from the KBS LTER and their respective satellite images were retrieved. Data preprocessing for the crop yield data is performed farm-wise, year-wise, 2x2 pixel-wise. Data preprocessing of the multi-



temporal satellite data is performed. This is done in order to build the dataset that is fit to be used in a deep learning model.

Two Regression Deep Learning models were implemented. The first one being the model with the input features comprising seven spectral bands namely: blue, red, green, near-infrared, short wave infrared 1, short wave infrared 2, and thermal infrared. The second one being the model with the input features comprising of Vegetation indices such as: NDVI, EVI2, SAVI, GRNDVI.

Further data augmentation is carried out for both the models to increase the data set and the accuracy of the model. This step includes image rotation, image flipping and cropping into 2x2 pixel windows. Annotation of crop yield corresponding to the seven spectral bands and four vegetation indices were done for all the farm images and for the 2x2 cluster images.

The dataset that has been created is segregated into training dataset, validation dataset, testing dataset for both the seven spectral bands and four vegetation indices. Further the input for the fine-tuned regression deep learning model 1 and 2 as shown in the flowchart is the training dataset along with validation dataset. The deep learning model gets trained based upon the training dataset. Then it is tested using the testing dataset. The crop yield of the individual farm is predicted using 2x2 farm yield. The accuracy of both the regression deep learning model 1 and 2 is compared and the results are analyzed.

### **3.2 STUDY AREA**

There is no readily available single dataset that can be downloaded and used for crop yield prediction in Michigan. The primary dataset used for predicting the crop yield in the state of Michigan was acquired from Kellogg Biological station. The crops (corn, soy, and wheat) are harvested annually using a combine machine equipped with GPS and precision agriculture software to allow detailed geo-referenced crop yield measurements with coincident GPS latitude and longitude data.

The raw data consisted of the geo-referenced grain flow rate and other parameters which are processed to calculate yield. This resulted in the processed dataset as shown in table 3.1. The dataset consisted of latitude and longitude points, along with the corresponding yield, harvest date and year, and moisture content. The year of harvest in the experimental farm is between the years 2001-2012 excluding the year 2002.

***Table 3.1 Dataset description***

<b>Variate</b>	<b>Description</b>	<b>Units</b>
longitude	longitude of the monitor at the time of reading	degree
latitude	latitude of the monitor at the time of reading	degree
crop yield	Yield as measured by sensor	bushelsPerAcre
moisture	gravimetric moisture of grain	%
Species	Type of crop planted	-
Year	Year the crop was harvested	-

The dataset consisted of 2,85,153 data points, which was split into 11 years according to the year of plantation and data that was scripted in excel. These 11 excel files were converted into the KML format. The intention was to find the number of farms under study for the time frame between 2001 to 2012(excluding the year 2002). KML is a file format used to display geographic data in an Earth browser such as Google Earth Pro. After analyzing the geographic data using the KML files, using Google Earth Pro for 11 years (2001-2012 (excluding the year 2002)) it was observed that the same 24 farms were under study as seen in Fig 3.2.



*Fig 3.2 Geographic data plotted in Google Earth Pro for the year 2001*

The four latitude and longitude corner data points for each farm were manually analyzed and tabulated using google earth pro with the help of geographic data as shown in Fig 3.2.

A farm has multiple data points based on the position and movement of the combine machine during harvest across the farm. Each data point representing a latitude and longitude points to a particular yield value which is based on the flow rate of the grain passing the sensor (crop flow as mentioned in the raw data). Each of the data points has a crop yield value that is represented in Bushels Per Acre. As a result, 264 yield data values were obtained corresponding to the 24 farms for all 11 years after the data preprocessing of the crop yield data as discussed in following sections.

### **3.2.1 CROPS IN STUDY AREA**

The crops under study in the KBS LTER farms are corn, soybean and wheat.

Corn being a warm season crop, it is preferred to be planted in the month of May when peak summer is expected. Plantation of corn in cool conditions results in imposing stress on the crop and affects the seedling health. Therefore, it is planted in the month of May.

Soybean is a warm season crop that is planted in the mid of May till the beginning of June. The land is prepared by plowing and laddering in order to sow the soybean seeds. Soybeans produces a better yield in warm and moist climate.

Wheat is a cold season crop. It is a thermally sensitive crop, whose plantation is done in the end of September until mid of October. It is harvested in the month of July of the upcoming year as it requires warm temperature during harvesting.

The plantation and harvest dates of the crops under study are summarized in the table 3.2.

***Table 3.2 Planting and Harvesting period of the crops:***

<b>Crop</b>	<b>Planting month</b>	<b>Harvesting month</b>	<b>Duration (in months)</b>
Corn	May – June	Oct – Nov	5
Wheat	Sept - Oct	July – Aug	8
Soybean	May - June	Oct - Nov	4

The Raw dataset collected from KBS LTER provided harvest date of the crops understudy for the years 2003 to 2012 and 2001 as shown in table 3.3

***Table 3.3 Year wise Harvest Dates of the crop***

<b>YEAR</b>	<b>CROP</b>	<b>HARVEST DATE</b>
2001	Wheat	16-07-2001
2003	Soybean	02-10-2003
2004	Wheat	12-07-2004

2005	Corn	10-10-2005
2006	Soybean	10-10-2006
2007	Wheat	09-07-2007
2008	Corn	20-10-2008
2009	Soybean	16-10-2009
2010	Wheat	14-07-2010
2011	Corn	14-11-2011
2012	Soybean	11-10-2012

### 3.3 RETRIEVAL OF SATELLITE DATA

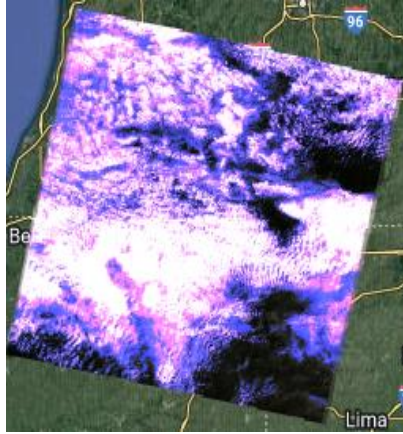
The approach of identifying and monitoring an area's physical features by measuring the reflected radiation from a distance is known as remote sensing. It generates a time series of spatially explicit data that can be used to characterize crop growth conditions and as input to crop models. These sensors capture data in the form of images, which may be manipulated, analyzed, and visualized using dedicated software. These images help in identifying minuscule details that aren't visible to human eyes.

In the study, we have chosen Landsat-5 and Landsat-7 satellite images (Level 2, Collection 2, Tier 1) with a spatial resolution of 30m. This dataset contains atmospherically corrected surface reflectance and land surface temperature derived from the data produced by the Landsat 5 TM sensor and Landsat 7 ETM+ sensor. The Landsat satellite sensors were designed to collect data in a different range of frequency bands across the electromagnetic spectrum. The Multispectral Scanner (MSS) present on Landsat 5 collected data in seven band ranges including a thermal and a shortwave infrared band with a spatial resolution of 30m for bands 1 to 7 excluding 6, whereas for band 6 the spatial resolution was resampled to 30m pixels from 120m pixels. Landsat 7's Enhanced Thematic Mapper Plus (ETM+) sensor images consist of eight spectral bands

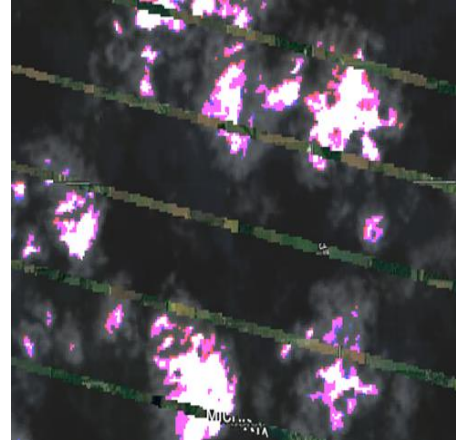
with a spatial resolution of 30 meters for Bands 1 to 7. The resolution for Band 8 which is panchromatic is 15 meters. Enhanced Thematic Mapper Plus (ETM+) sensor of Landsat 7 provides images that consist of eight spectral bands with a spatial resolution of 30 meters for Bands 1 to 7. The spatial resolution of the panchromatic band(Band 8) is 15 m.

The satellite images were collected for all the 24 farms each of the 11 years from Google Earth engine which is a cloud-based geospatial analysis platform that helps to visualize and analyze satellite images and process it according to our requirement using the tools and in-built commands. A JavaScript was written in GEE, to retrieve each farm as a GeoTIFF file. A GeoTIFF is a public-domain metadata standard for incorporating georeferencing information into image files. GeoTIFF embeds geospatial metadata into image files like aerial photography, satellite imaging, and digitized maps, allowing them to be used in GIS applications. This metadata standard allows the geographic data to be associated with the image data whereas a TIFF file standard allows the metadata as well as the image data to be encoded in the same file. This is the reason the TIFF file is used for storing raster images.

Collection of the satellite images from both Landsat 5 and 7 was performed in order to avoid the irregularities present in the image such as stripes as shown in Fig 3.4 and cloud cover of more than 10% as shown in Fig 3.3. Therefore, the appropriate satellite data images with cloud cover lesser than 10% were collected before the harvest dates without any of the above-mentioned problems. To maintain uniformity between the satellite images collected from Landsat 5 and 7, only 7 spectral bands of spatial resolution 30m were collected from the Google earth engine.



**Fig 3.3** Cloud cover of 93% (Landsat 5)



**Fig 3.4** Stripes in image (Landsat 7)

In order to fulfill the objective of predicting farm's yield right before harvest, satellite images of farms were chosen before the date of harvest in each year from 2001 to 2012. The total acquired dataset contained 252 images in total of the farms under study, over 11 years. All the farm images have pixel dimensions of either 4x4 or 4x5 or 5x4 or 5x5. In Fig 3.5, the "natural color" band combinations are observed. The visible bands are used in these band combinations.



**Fig 3.5** Farms Visualization in RGB Format

### 3.3.1 HISTORICAL DATA RETRIEVAL

The historical data would help to analyze the change in features of the farm, feature correlation, and temporal analysis from plantation to harvest period. The GeoTIFF image was extracted after the plantation of the crop. This helped the model learn how the farm has evolved from plantation to harvest.

The images were extracted based on the duration between the plantation and harvest dates of the crop. The historical image as shown in table 3.4 is taken in the initial periods of plantation whereas the other image was retrieved before the harvest. For every farm, two images are retrieved, one from the second half of the duration between plantation and harvest, and the other image from first half of the duration between plantation and harvest. The duration between the plantation and harvest date for wheat is 8 months. One image was retrieved in the last four months before harvest, as shown in Fig 3.7, whereas the other image was retrieved within the first four months after plantation, shown in Fig 3.6. Similarly for soybean, one image was retrieved in the last two months before harvest whereas the other image was retrieved within the first two months after plantation. For corn, one image was retrieved in the last 2.5 months before harvest whereas the other image was retrieved within the first 2.5 months after plantation.



**Fig 3.6** After plantation  
for the year 2001 (Wheat)



**Fig 3.7** Before harvest  
for the year 2001 (Wheat)



*Table 3.4 Landsat images retrieval dates*

<b>Crop</b>	<b>Year</b>	<b>Date of Image after plantation</b>	<b>Date of Image before harvest</b>	<b>Duration (in months)</b>
Wheat	2001	January 8	July 11	8
Soybean	2003	June 23	September 11	4
Wheat	2004	November 30 (2003)	July 11	8
Corn	2005	July 30	September 8	5
Soybean	2006	August 2	October 5	4
Wheat	2007	November 6 (2006)	June 18	8
Corn	2008	July 6	October 10	5
Soybean	2009	July 9	September 27	4
Wheat	2010	November 14 (2009)	June 10	8
Corn	2011	September 1	November 4	5
Soybean	2012	July 9	September 11	4

### **3.4 COMPUTATION OF CROP YIELD**

With the processed combine machine crop yield dataset that was collected as discussed in section 3.2, the farm-level crop yield was computed with the help of a python script that was coded to calculate the farm-wise yield value for each of the 11 years in terms of bushels per acre. The algorithm used to calculate the yield for a particular farm was to take an average of the yield values of all the data points within the boundary of each farm. In section 3.7, the computation of the crop yield of all the 4 pixel windows of the farms will be performed.

#### **3.4.1 LOG TRANSFORMATION TECHNIQUE**

Since the crop yield values of the farms have a wide value range, logarithmic transform technique was applied to it in order to facilitate the ease of learning of the regression

deep learning model. This is due to the reduced range of crop yield values after applying Log transform.

### 3.5 EXTRACTING BAND VALUES

The 7 Band values of the acquired GeoTIFF files were extracted using the rasterio library. Rasterio reads and writes these formats and provides a Python API based on NumPy N-dimensional arrays. Each farm's dimension was approximately 105m x 87m. With the spatial resolution of the Landsat satellite being 30m, each farm's image had around 16 to 25 pixels. With 252 images in the dataset, the number of features for each image ranged from 112 to 175. Say, for instance, a farm is having 16 pixels, a total number of input features would be 112 (7 bands \* 16 pixels). This dataset is too small to perform any deep learning methods on it since there exists only 252 images under study. The following bands were collected as shown in table 3.5

*Table 3.5 Multispectral Bands*

<b>Band</b>	<b>Description</b>	<b>Wavelength</b>	<b>Spatial Resolution</b>
SR_B1	Blue	0.45-0.52 $\mu$ m	30m
SR_B2	Green	0.52-0.60 $\mu$ m	30m
SR_B3	Red	0.63-0.69 $\mu$ m	30m
SR_B4	Near-Infrared	0.77-0.90 $\mu$ m	30m
SR_B5	ShortWave Infrared 1	1.55-1.75 $\mu$ m	30m
SR_B7	ShortWave Infrared 2	2.08-2.35 $\mu$ m	30m
ST_B6	Thermal Infrared	10.40-12.50 $\mu$ m	30m

The significance of the spectral bands is explained below:

Band 1: Blue band distinguishes between soil and vegetation. Also, distinguishes deciduous from coniferous vegetation.

Band 2: Green band detects peak vegetation areas and is useful for the measure of the

increase in plant growth.

Band 3: Red band distinguishes between vegetation slopes. Also, it can draw a comparison between cropland with standing crops and cropland with stubble.

Band 4: Near-infrared band distinguishes between barren land and cropland since it emphasizes biomass content.

Band 5: Short wave infrared -1 band distinguishes the moisture content of the soil and vegetation. It also separates forest lands, croplands, and water bodies.

Band 6: Thermal infrared band estimates the soil moisture content

Band 7: Short wave infrared -2 band distinguishes between water and rocks, hence useful for mapping hydrothermally altered rocks. Also, it can distinguish between land and water.

As observed, all 7 bands can contribute to crop yield prediction since the main factors involved in the prediction of yield is vegetation, moisture content, to identify barren / cropland/forests, etc. Hence, it is concluded to have these seven bands as the input features.

### **3.6 EXTRACTING VEGETATION INDICES**

Vegetation index is a method of feature engineering in which various bands are combined to form a compact and reasonable single feature. These are spectral transformations of satellite bands that help to analyze the vegetative properties. It results in accurate temporal and spatial comparisons of the photosynthetic activity and variations in canopy structure.

By incorporating vegetation indices with a deep learning model or a software the current limitations in the agricultural sector can be resolved. This is achieved by performing analysis on the satellite images that is used to examine the climatic trends, remotely estimate the water content in soils, classify vegetation, monitoring droughts, scheduling crop irrigation, crop management and accessing any changes in biodiversity. The vegetation index is further classified into Multispectral vegetation indices and Hyperspectral vegetation indices. In this study, the main focus is on Multispectral vegetation indices.

The indices used in this study are:

### 1. NDVI (Normalized Difference Vegetation Index)

Health of the vegetation, differentiating between vegetation types, detecting unusual and abnormal changes in the process of plantation to harvest, indicating the active biomass, are few factors that are indicated by the NDVI (Normalized Difference Vegetation Index). NDVI values can range from -1.0 to +1.0. It is the most important vegetation index that is used in agriculture applications. The mathematical formula for NDVI calculator, is derived from the BAND 4 (Near-infrared) and BAND 3 (Red) respectively:

$$\frac{NIR - RED}{NIR + RED} \text{ or } \frac{band4 - band3}{band4 + band3} \quad (1)$$

### 2. SAVI (Soil Adjusted Vegetation index)

SAVI vegetation index helps in analysis of crops in regions where vegetation cover is low. An adjusted factor L has been introduced to NDVI which eliminates the soil noise effects that were previously present in NDVI. So, this is said to be a soil adjusted vegetation index. The factor L value depends on green vegetation density. Its value varies between -1 to +1. L=0.5 is chosen to adjust land cover. The mathematical formula used for Landsat-5 and Landsat-7 is:

$$\frac{NIR - RED}{NIR + RED + L} \times (1 + L) \text{ or } \frac{band4 - band3}{band4 + band3 + 0.5} \times (1.5) \quad (2)$$

### 3. EVI-2(Two band Enhanced Vegetation Index):

EVI is a vegetation index that indicates and quantifies the greenness of a region. It is similar to NDVI, but it takes in an additional feature to eliminate the canopy background noise, and compensates for noise due to atmospheric conditions. EVI-2 indicates the chlorophyll content with the least topographic effects. EVI-2 eliminates BAND-1 (blue) since it has a poor ratio of S/N (signal to noise ratio). The range of EVI-2 varies between 0.2-0.8. The mathematical formula to calculate EVI-2 for Landsat-5 and Landsat-7 is:

$$2.5 \times \frac{NIR - RED}{((NIR) + (C1 * RED) + L)} \text{ or } 2.5 \times \frac{NIR - RED}{((NIR) + (2.4 * RED) + 1)} \quad (3)$$

The value for C1 is chosen to be 2.4. This is used to compensate for the aerosol scattering in the atmosphere whereas L is used to adjust the noise factors.

#### 4. GRNDVI (Green Red NDVI):

Modified NDVI with green and Red bands in GRNDVI. It helps to analyze the various aspects of crop like chlorophyll content, vegetation at maturity of the crop. The range of GRNDVI varies from -1.0 to +1.0. The mathematical formula used to calculate GRNDVI for Landsat-5 and Landsat-7 is:

$$\frac{NIR - (GREEN + RED)}{NIR + (GREEN + RED)} \text{ or } \frac{band4 - (band2 + band3)}{band4 + (band2 + band3)} \quad (4)$$

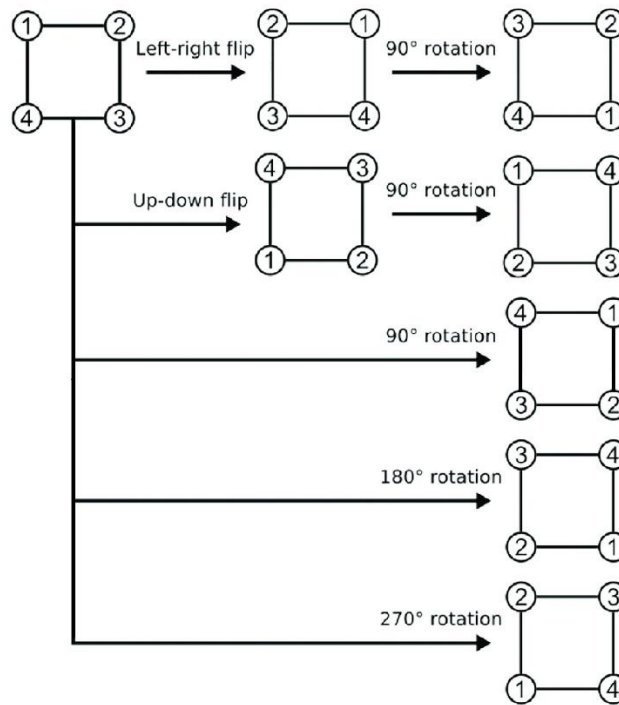
### 3.7 METHODS TO INCREASE DATASET

To increase the dataset, each farm was cropped out into a cluster of adjacent pixels using a 4-pixel window (2x2 pixels). The adjacent pixels are chosen to implement this window in order to incorporate the covariance factor in the image dataset. This created a possibility of acquiring 9 to 16 different images (each of size 2x2 pixels) based on the number of pixels of the farm from the downloaded GeoTIFF files. A python script was coded to calculate each of the pixel's latitude and longitude endpoints for all of the 2x2 farm images. This is calculated using the gdal library. In order to find out the yield corresponding to each 2x2 image, the pixel-wise yield was computed with the help of the pixel-wise endpoint latitude and longitude coordinates using the python script. By making use of pixel-wise yield, the yield corresponding to each of the 2x2 pixel images were calculated. This technique was performed on both the historical dataset as well the original dataset. This resulted in an increased dataset with approximately 2,884 images with their respective yield values.

### 3.7.1 AUGMENTATION

Augmentation is a technique used to increase the size of the dataset by performing simple operations on the existing dataset such as rotation, shearing, zooming, cropping, flipping, and changing the brightness level. In this study, two methods of augmentation namely: rotation and flipping were performed on each of the 2x2 pixel images. As shown in Fig 3.8, about 8 augmented images are being produced for each 2x2 pixel image by performing rotation and flipping. Hence, it has increased the dataset 8 times compared to the original dataset, resulting in 23,072 images.

It is the process of artificially expanding the available dataset for training a deep learning model. Techniques such as rotation and flipping were used.



*Fig 3.8 Augmentation*

### 3.8 ANNOTATION

The augmented datasets of 7-spectral bands and four vegetation indices are the two training X datasets for the regression deep learning model 1 and regression deep learning

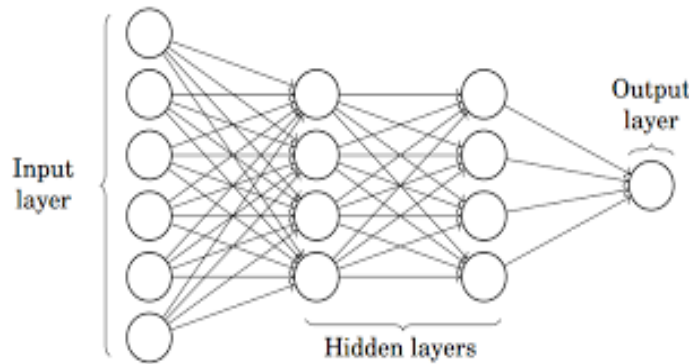
model 2 respectively as discussed in section 3.1 (workflow). The training X datasets are labeled with the corresponding crop yield data of the 4-pixel window of the farm which was calculated in the data preprocessing of the crop yield data.

### 3.9 DEEP LEARNING:

Deep Learning modern technology is a subset of machine learning which again is a subset of Artificial intelligence. Deep Learning architecture was developed to mimic the behavior of the human brain. This allows it to learn the intricate features by analyzing large amounts of data with the given logic structure.

There are various artificial intelligence (AI) applications driven by deep learning to improve various sectors such as automation, healthcare, fraud detection, Entertainment, Visual Recognition, and Natural Language Processing. To implement the above, deep learning makes use of the multi-layered structure of algorithms called neural networks.

### 3.10 NEURAL NETWORK:



*Fig 3.9 Neural Network*

A neural network consists of three stages/layers namely the input layer, hidden layers, and output layer. 1. Input layer, consists of nodes that take in raw information from the outside world and propagate it to the hidden layer. 2. Hidden layers consist of neurons, where all computation and feature extraction/learning of the model is processed. 3.

Output layer provides the output of the model, or conclusion obtained after processing the input.

Neural networks can be used for variety of tasks, including clustering, classification, and regression. The layers can independently learn an underlying representation of the raw input. In subsequent layers of artificial neural networks, a more abstract and condensed representation of the raw data is created. The result is then derived using this compressed form of the input data. This implies that the feature extraction stage is already included in the artificial neural network process. In the training process, the above-mentioned process is optimized by the neural network in order to create the best possible abstract representation of the input data.

An artificial neural network is made up of a series of interconnected units or nodes or neurons. These artificial neurons are similar to the biological human brain neurons. A neuron is nothing but a graphical representation of a numeric value. The connections of biological neurons are called axons. The axons in the artificial neural networks are termed “weights” which are also numeric values. These weights change when the model learns from the data. This process of the model learning to find the weights for the given dataset is called “training”.

### **General working:**

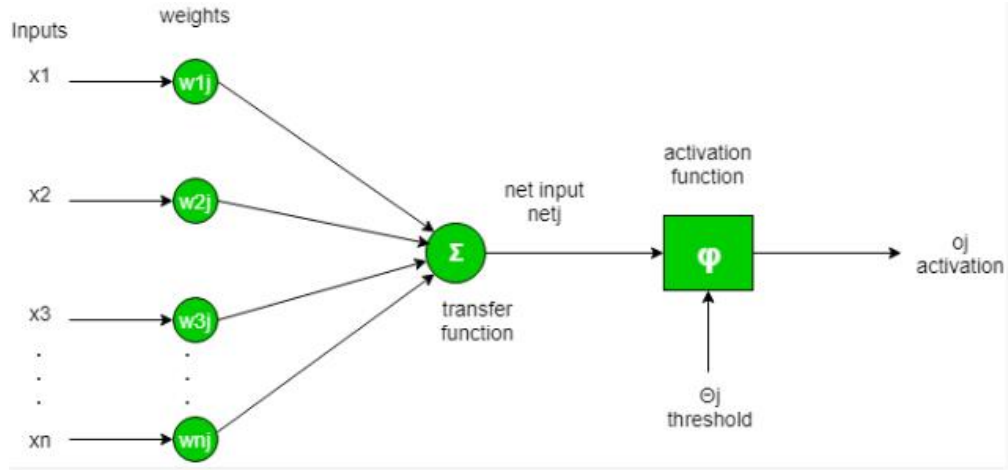
Initially, inputs  $X_1, X_2, \dots, X_n$  are passed along with their assigned weights to the hidden layer (which consists of multiple neurons). All the inputs are connected to each of the neurons, and in the hidden layer, computation takes place.

Consider a single artificial neuron in the hidden layer, where the inputs are  $X_1, X_2, \dots, X_n$ . All inputs are connected to the neuron, where the processing and computations take place.



### Computation in neuron:

A neuron performs computation by adding the weighted sum of all inputs and the bias value. Computation steps in a neuron can be implemented in two steps as mentioned below in Fig 3.10



**Fig 3.10** Computation in neuron

Weights are assigned to the inputs of the neuron based on the importance/ priority of a particular input. It is found by finding the gradient of an input variable.

$$Z_1 = W_1 \times X_1 + W_2 \times X_2 + \dots + W_n \times X_n \quad (5)$$

A bias value is added to  $Z_1$  to help fit the model in a better way.

$$Z_1 = W_1 \times X_1 + W_2 \times X_2 + \dots + W_n \times X_n + b \quad (6)$$

This is the first step in the computation process which can mathematically be represented by:

$$\sum_{i=1}^n X_i \times W_i \times b \quad (7)$$

Here,  $W_i$  represents weight,  $X_i$  represents input features,  $b$  represents the bias value and  $n$  represents the number of input features.

1. The value of the computation done in step one can be in the range of  $-\infty$  to  $+\infty$ . The activation function is a threshold function that bounds the value and helps to decide if a neuron needs to be activated or fired. There are various types of activation functions such as ReLU (Rectified Linear Unit), Step function, and Sigmoid function.
2. The above process is repeated till the output from the final layer's neuron is received. After this prediction, the output should be compared with the real label of the input. The loss function defines the difference between the predicted output and the real label. Finding the global minima of the loss function will help to increase the accuracy of the model since the difference between the predicted and the real label of the input decreases. There are various loss functions such as quadratic loss, cross-entropy loss, and mean squared error loss.
3. The mathematical procedure to minimize the loss function is called gradient descent. The stochastic gradient descent technique is used to train deep learning neural networks. Stochastic gradient descent is an optimization algorithm that calculates the gradient error for the model using the training dataset and changes the weights of the model using the back-propagation technique.
4. The number of times the training data is shown to the neural network is called epochs. Overfitting occurs when the epoch is increased beyond a point where validation accuracy decreases even when the training accuracy is increasing.
5. Batch size is the number of subsamples that are given to the model where hyperparameter tuning happens.

### **3.11 ANN MODEL - SPECTRAL BANDS**

Artificial regression deep learning neural network model was trained using the pool-seven spectral band dataset. A sequential model from the keras TensorFlow library was used to build a deep learning model. The first layer has input dimensions of 56 since the number of features of a single data point is calculated by two multi-temporal farm images of 4-pixel size for the 7-spectral bands ( $2 \times 4 \times 7 = 56$ ). The output layer has a single neuron since it is a regression deep learning model. The dataset was standardized using the sklearn StandardScalar library. It was then divided into training, validation, and testing dataset in the ratio of 80:10:10 such that the testing dataset has all of the four-pixel

window combinations which when put together gives the total pixels of the farm image for the selected random pool farms. Using the grid search method while training the model, tuning of hyperparameters such as number of neurons, number of dense hidden layers, number of epochs, and batch size was performed with the help of the validation dataset. The tuning was based on minimizing the Mean squared error. Based on the tuning, the architecture of the model was with 10 dense layers with the neurons sequence as 2048, 1024, 512, 256, 128, 64, 32, 16, 8, 1. The batch size and number of epochs was 200 and 450 respectively. The Adam optimizer was used as the stochastic gradient descent method for the loss function which has a learning rate of 0.001. The ReLU activation function was used for all of the layers except for the output layer where the linear activation function was used since ReLU is the most popular activation function that is used that doesn't have the issue of vanishing gradient. The kernel initializer was set to be a normal distribution. The loss vs epoch graph was plotted for both the training and validation dataset. The error of both the training and validation graphs was decreasing until the fine-tuned epoch number which showed that the model is not overfitting with the dataset.

The model output after the inverse transform of the standardization predicts the crop yield for the 4-pixel window of the total farm as a logarithm value of the original crop yield as discussed in section 3.4.1 Anti-log transform was applied to the predicted output to obtain the crop yield in terms of bushelsPerAcre. RMSE, MAPE, and R squared score was calculated for evaluating the performance of the deep learning model.

The full farm yield is calculated using the following algorithm:

1. Predict the 2x2 image crop yield of the testing data using the proposed regression deep learning model.
2. The four-pixel windows of the farm have the yield in terms of bushelsPerAcre. Thus, in order to calculate the average crop yield of the total farm the weighted average of all of the 4-pixel windows' crop yield has to be found.
3. Some of the four-pixel windows of the farm are situated at the edges of the farm with the pixel dimensions of 4x5, 5x4, and 5x5 are prone to overlap with the other 4-pixel

windows of the same farm. These four-pixel windows are given the weights of non-overlapping pixel size.

**Table 3.6** Hyperparameter tuning – Spectral band model:

Batch size \ Epoch	50	75	100	150	200	250
100	89.02	90.46	92.10	89.32	89.59	89.62
200	92.23	94.14	94.97	94.78	94.60	91.40
300	89.01	95.42	93.61	95.47	96.23	95.81
350	94.732	94.83	94.02	95.75	95.60	96.02
400	94.99	93.37	96.00	94.42	94.47	96.19
450	95.47	95.28	95.52	95.55	97.35	95.62
500	92.69	96.37	95.79	96.01	96.31	95.89

### 3.12 ANN MODEL - VEGETATION INDICES

Artificial regression deep learning neural network model was trained using the pool-four vegetation indices dataset. A sequential model from the keras TensorFlow library was used to build a deep learning model. The first layer has input dimensions of 32 since the number of features of a single data point is calculated by two multi-temporal farm images of 4-pixel size for the 4-vegetation indices (NDVI, GRNDVI, EVI2, SAVI) ( $2 \times 4 \times 4 = 32$ ). The output layer has a single neuron since it is a regression deep learning model. The dataset was standardized using the sklearn StandardScalar library. It was then divided into training, validation, and testing in the ratio of 80:10:10 such that the testing dataset has all of the four-pixel window combinations which when put together gives the total pixels of the farm image for the selected random pool farms. Using the grid search method while training the model, tuning of hyperparameters such as number of neurons, number of dense hidden layers, number of epochs, and batch size was performed with the help of the validation dataset. The tuning was based on minimizing the Mean squared

error. Based on the tuning, the architecture of the model was with 10 dense layers with the neurons sequence as 2048, 1024, 512, 256, 128, 64, 32, 16, 8, 1. The batch size and number of epochs was 100 and 600 respectively. The Adam optimizer was used as the stochastic gradient descent method for the loss function which has a learning rate of 0.001. The ReLU activation function was used for all of the layers except for the output layer where the linear activation function was used since ReLU is the most popular activation function that is used that doesn't have the issue of vanishing gradient. The kernel initializer was set to be a normal distribution. The loss vs epoch graph was plotted for both the training and validation dataset. The error of both the training and validation graphs stayed decreasing until the fine-tuned epoch number which showed that the model is not overfitting with the dataset.

The model output after the inverse transform of the standardization predicts the crop yield for the 4-pixel window of the total farm as a logarithm value of the original crop yield as discussed in section 3.4.1. Anti-log transform was applied to the predicted output to obtain the crop yield in terms of bushelsPerAcre. RMSE, MAPE, and R squared score was calculated for evaluating the performance of the deep learning model. The full farm yield is calculated using the same algorithm as discussed in section 3.11.

**Table 3.7** Hyperparameter tuning – Vegetation Indices model:

Batch size \ Epoch	50	75	100	200	250
300	92.11	92.52	93.44	93.71	90.42
350	92.29	95.22	94.88	92.98	87.86
400	96.00	96.12	95.24	95.94	95.33
450	92.91	95.39	93.82	87.27	92.85
500	90.41	95.29	94.59	95.73	96.57
600	96.06	94.18	96.77	96.62	96.24
700	96.57	92.89	92.50	95.74	94.80

## **CHAPTER 4: RESULTS & DISCUSSIONS**

## **4.1 SOFTWARE PLATFORMS**

Google code earth engine platform was used to download the Landsat-5 and Landsat-7 satellite dataset of the study area. The python script for the entire proposed model was coded in Google colab's jupyter notebooks. The entire dataset for the proposed model was stored in Google Drive which is mountable with Google Colab. Google Colab is a cloud-based network, and it allocates virtual machine with 2v CPU (Intel(R) Xeon(R) CPU @ 2.30GHz), 13GB RAM, 40GB Hard disk drive and Nvidia Tesla K80/T4 GPU with maximum execution time of 12 hours.

## **4.2 RESULTS**

Two regression deep learning models were built to predict the crop yield at the farm level. The study area chosen was the experimental farms in Michigan which have 24 farms in total for the span of 11 years. One of the deep learning models was trained using the spectral bands' dataset as features while the other deep learning model with the four vegetation indices dataset as features.

### **4.2.1 DATASET PREPARATION**

The dataset contains a total of 252 Landsat satellite farm images over the span of 11 years (2001-2012 excluding 2002). The corresponding crop yield data of the farm were calculated from the average crop yield data which was observed from the combine machine. To increase the correlation of the crop growth at different time period, a historical image is retrieved such that one of the images lies in the first half of the plantation period while the other image lies in the second half of the plantation period. Hence, the original dataset consists of 252 farm images, which is very less number of data points for training the deep learning model. In order to increase the data points, the images were cropped in a sliding four-pixel window and their corresponding 4-pixel crop yield was also computed using the farms' geo-referenced co-ordinates. This increased the dataset to 2884 data points as shown in table 4.1

**Table 4.1** *Size of dataset before Augmentation*

<b>DATASET SIZE</b>	<b>SIZE</b>
Original Dataset	252
2x2 pixel images	2,884

Using all the 4-pixel images, two separate datasets were created. One of the dataset is the 7- spectral bands (Band 1 (blue), Band 2(green), Band 3 (red) , Band 4 (near infrared), Band 5 (shortwave infrared 1) , Band 6 (thermal infrared)) and the other being the 4-vegetation indices (NDVI, GRNDVI, EVI-2, SAVI). Augmentation of the two datasets were performed by rotating and flipping the respective data points. This resulted in an increased dataset of 23,072 data points as shown in table 4.2

**Table 4.2** *Size of dataset after augmentation*

<b>DATASET</b>	<b>SIZE</b>
2x2 pixel images	2,884
With augmented 2x2 pixel images	23,072

There are 2 datasets created, one being the 7 spectral band and the other being the 4 vegetation indices dataset. The feature size of the 7 spectral band dataset for a single 2x2 pixel image along with historical data is 56 (2 x 4 x 7). Similarly, the feature size of the vegetative indices dataset for a 2x2 pixel image is 32 (2 x 4 x 4)

**Table 4.3** *Feature size*

<b>DATA SET</b>	<b>FEATURE SIZE</b>
7 Spectral Band	56
4 Vegetation Indices	32



### 4.2.2 PERFORMANCE METRICS

Error metrics are used to quantify the error produced by the model. For the analysis, three error metrics are used namely RMSE, R2 score, and MAPE.

1. RMSE (Root mean squared Error)

$$RMSE = \sqrt{\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N}} \quad (8)$$

2. MAPE (Mean Absolute percentage error)

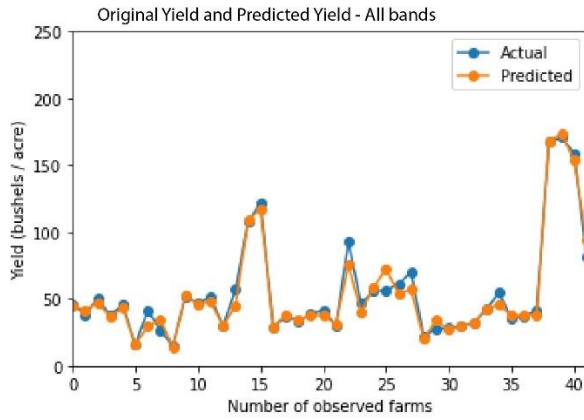
$$MAPE = \frac{1}{N} \sum_{i=1}^N \left| \frac{y_i - \hat{y}_i}{y_i} \right| \quad (9)$$

Here, N represents the number of data points,  $y_i$  represents the  $i^{th}$  actual value and corresponding  $\hat{y}_i$  represents the predicted value.  $\mu$  is the mean of actual values ( $y$ ).

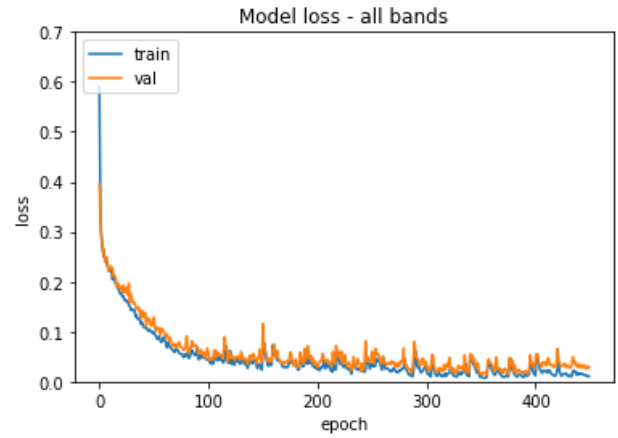
3. Accuracy

$$Accuracy = (100 - MAPE) \quad (10)$$

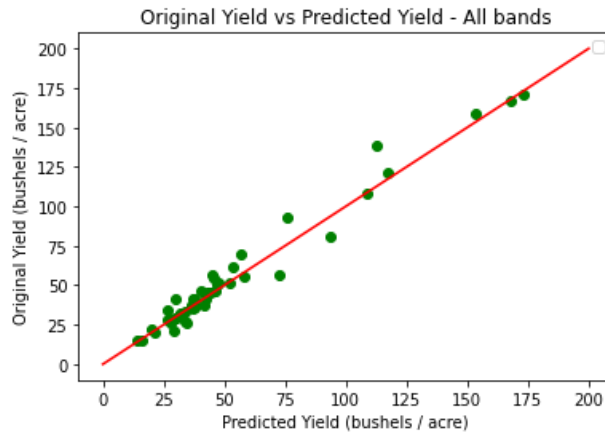
### 4.2.3 ANALYSIS OF POOL SPECTRAL BANDS



**Fig 4.1** Original Yield and Predicted Yield - All bands



**Fig 4.2** Model Loss - All bands



**Fig 4.3** Original Yield vs Predicted Yield - All bands

**Table 4.4** Performance of pool spectral band model

RMSE	MAPE	Full Farm Accuracy
7	5.8	94.2%

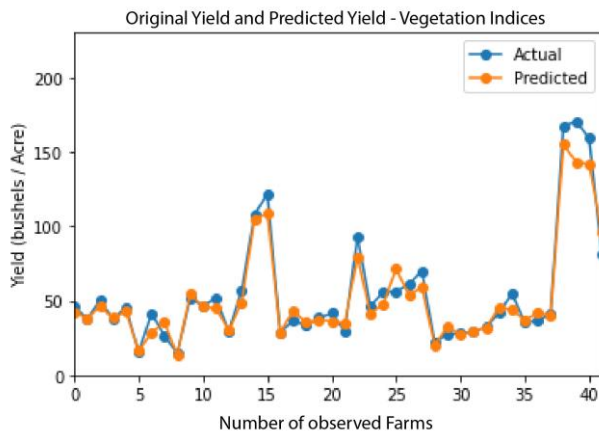
The loss vs epoch graph was plotted for the pool 7 spectral bands dataset for both the training and validation dataset as shown in Fig 4.2. The error of both the training and

validation graphs stayed decreasing until the fine-tuned epoch number which is 450 shows that the model is not overfitting with the dataset.

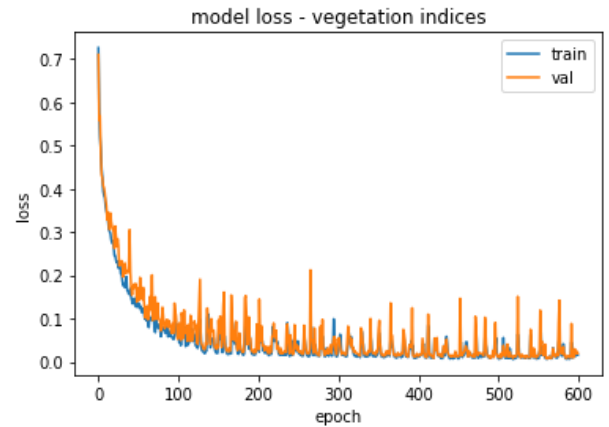
The original and the predicted yield are plotted for all the testing pool farms in Fig 4.1 which visualizes the accuracy of the predicted total farm yield. On the other hand, Fig 4.3 plots the scatter of original vs the predicted crop yields for all of the testing farms.

Table 4.4 has all the error calculations which are useful for analyzing the performance of the pool 7 spectral bands regression deep learning model. The RMSE of 7 for the crop yield dataset of range 0 to 200 is a very good score. The R squared score of 0.96 is also considered a good score which shows a high correlation of the predicted crop yield data with the original crop yield data. The mean absolute error (MAPE) and the accuracy of the model are 5.8 and 94.2 percent respectively. This shows the mean of the difference of error of the predicted crop yield is less.

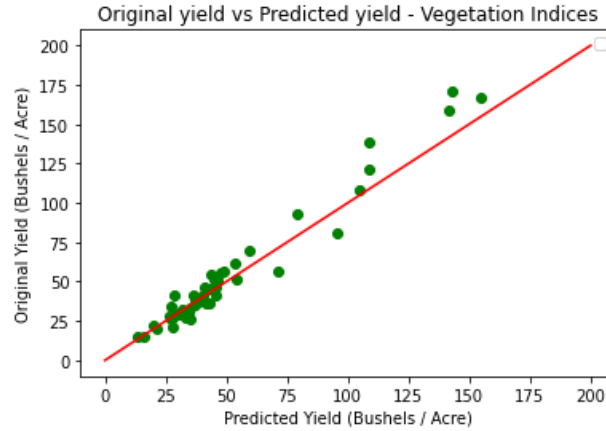
#### 4.2.4 ANALYSIS OF POOL-VEGETATION INDICES



**Fig 4.4** Original Yield and Predicted Yield - Vegetation Indices



**Fig 4.5** Model Loss -Vegetation Indices



**Fig 4.6** Original Yield vs Predicted Yield-Vegetation Indices

**Table 4.5** Performance of pool Vegetation Indices model

RMSE	MAPE	Full Farm Accuracy
8.2	11.2	88.8%

The loss vs epoch graph was plotted for the pool 4 vegetation indices dataset for both the training and validation dataset as shown in Fig 4.5. The error of both the training and validation graphs stayed decreasing until the fine-tuned epoch number which is 600 shows that the model is not overfitting with the dataset.

The original and the predicted yield are plotted for all the testing pool farms in fig 4.4 which visualizes the accuracy of the predicted total farm yield. On the other hand, Fig 4.6 plots the scatter of original vs the predicted crop yields for all of the testing farms.

Table 4.5 has all the error calculations which are useful for analyzing the performance of the pool 4 vegetation indices regression deep learning model. The RMSE of 8.2 for the crop yield dataset of range 0 to 200 is a very good score. The R squared score of 0.95 is also considered a good score which shows a high correlation of the predicted crop yield data with the original crop yield data. The mean absolute error (MAPE) and the accuracy of the model are 11.2 and 88.8 percent respectively. This shows the mean of the difference of error of the predicted crop yield is less.

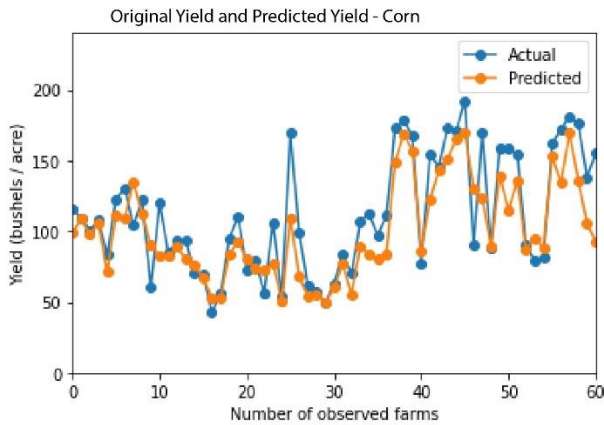
As observed in section 4.2.3 and 4.2.4, models trained using the 7-spectral produced an

accuracy of 94.3% whereas the model trained with vegetation indices as input features produced an accuracy of 88.8%.

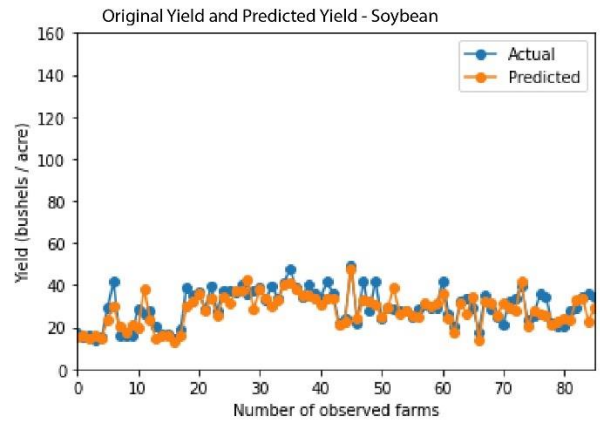
The seven spectral bands contribute unique features and hence the model can identify patterns of various features while testing. On the other hand, for the second model the vegetation indices used as input features for training the model were NDVI, GRNDVI, SAVI, EVI-2 which have great importance in agriculture applications. These indices are results of the computation of the three spectral band values as discussed in chapter 3. So, the model misses out on learning the patterns due to the other four bands since only vegetation indices were given as input features. Hence, the accuracy of the model-2 is less compared to that of the Pool-band model. (model-1)

So, further analysis of year-wise yield prediction and crop wise yield prediction was performed with input features to be the 7-spectral bands only.

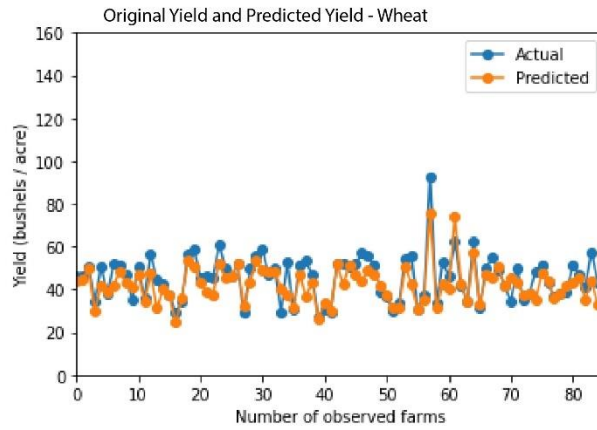
#### 4.2.5 CROP WISE ANALYSIS:



**Fig 4.7** Crop yield and Predicted yield for Corn



**Fig 4.8** Crop yield and Predicted yield for Soybean

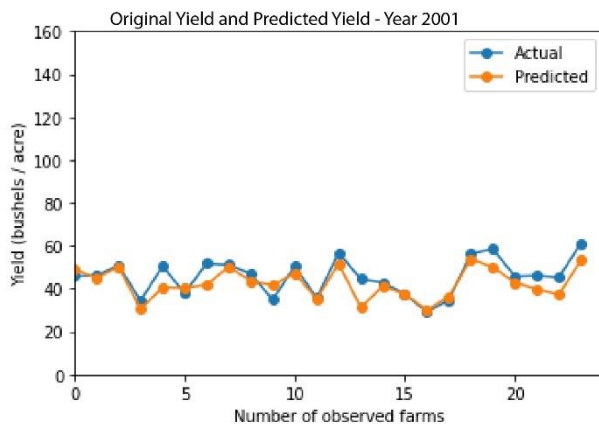


**Fig 4.9** Crop yield and Predicted yield for Wheat

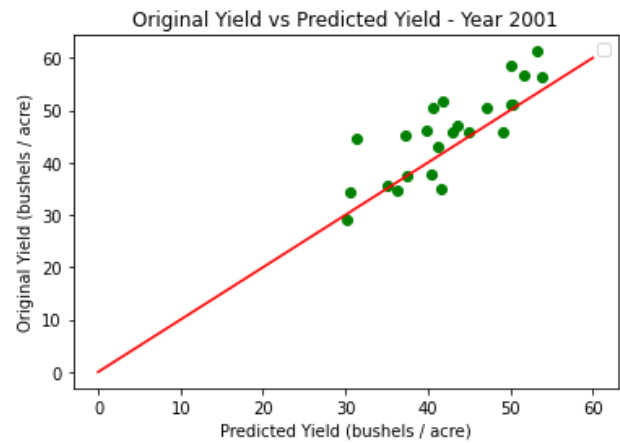
**Table 4.6** Crop-wise analysis

Crop	RMSE	MAPE	Full-Farm Accuracy
Wheat	6.58	10.27	89.7
Soybean	4.94	12.17	87.8
Corn	22.9	15.26	84.74

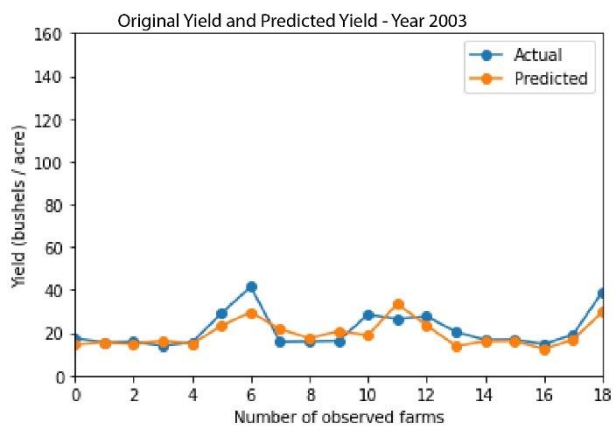
#### 4.2.6 YEAR-WISE ANALYSIS



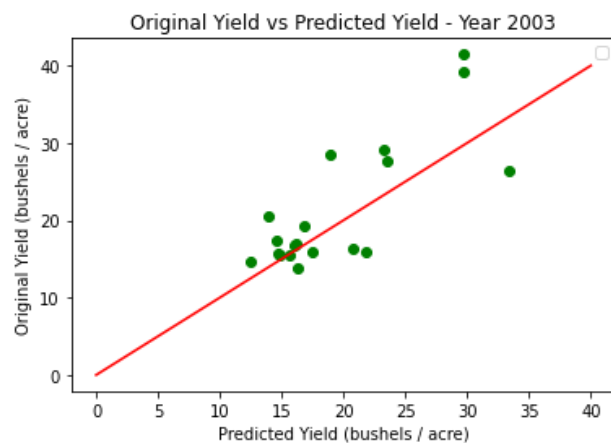
**Fig 4.10** Original yield and Predicted yield for the year 2001



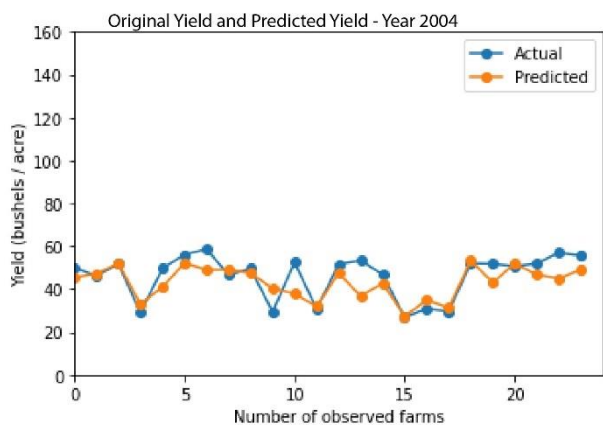
**Fig 4.11** Original yield vs Predicted yield for the year 2001



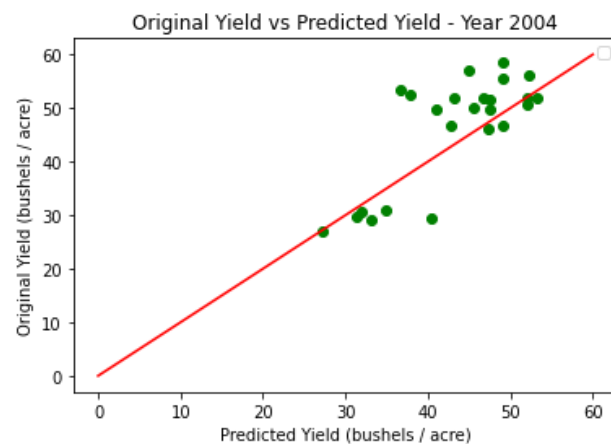
**Fig 4.12** Original yield and Predicted yield for the year 2003



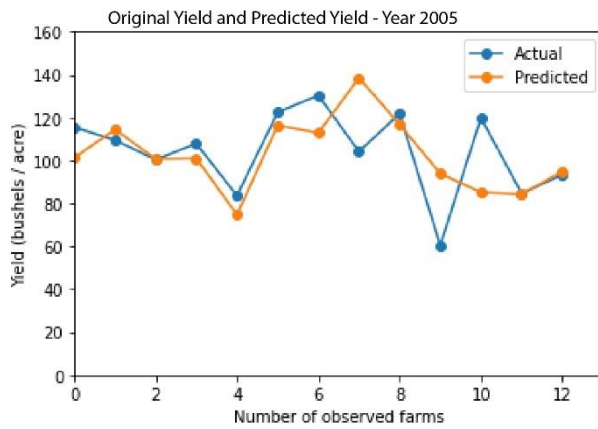
**Fig 4.13** Original yield vs Predicted yield for the year 2003



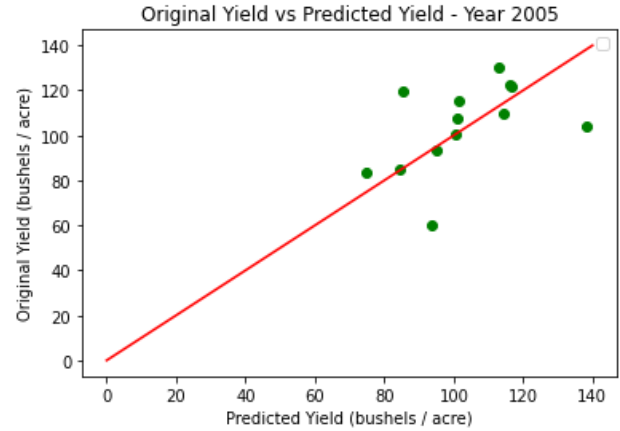
**Fig 4.14** Original yield and Predicted yield for the year 2004



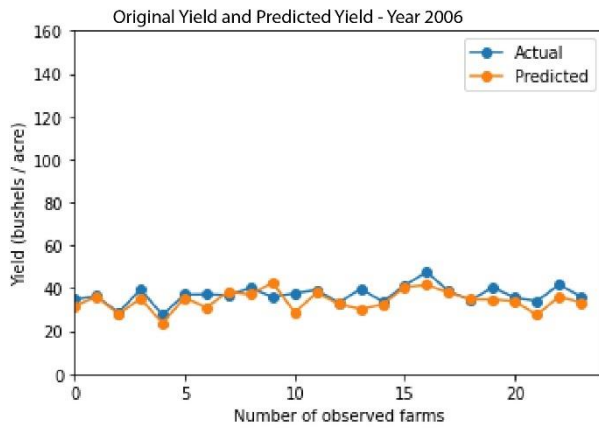
**Fig 4.15** Original yield vs Predicted yield for the year 2004



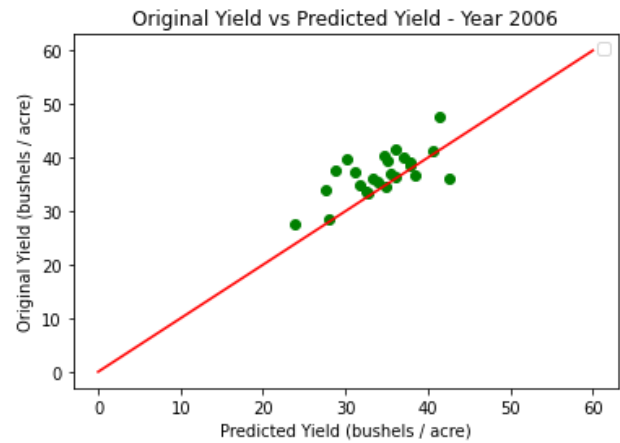
**Fig 4.16** Original yield and Predicted yield for the year 2005



**Fig 4.17** Original yield vs Predicted yield for the year 2005

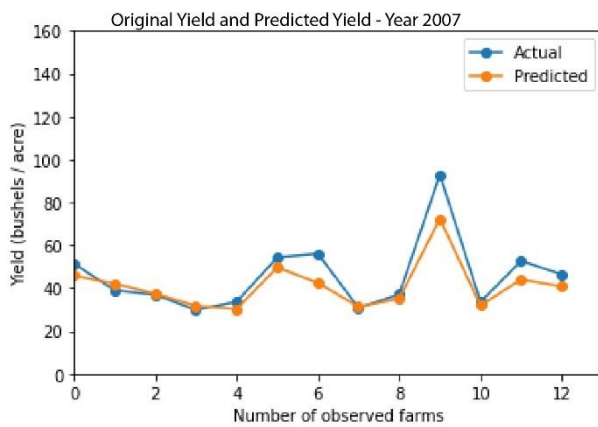


**Fig 4.18** Original yield and Predicted yield for the year 2006

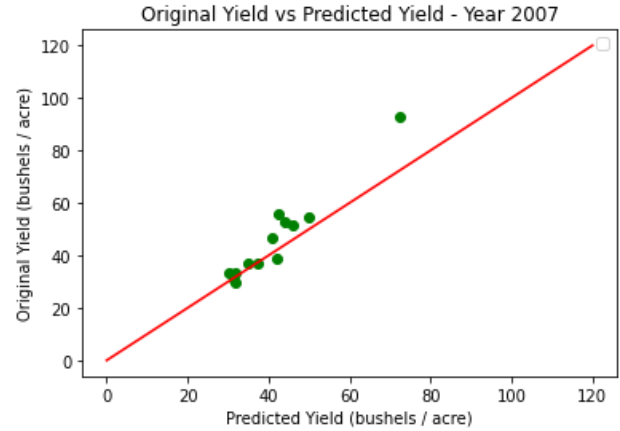


**Fig 4.19** Original yield vs Predicted yield for the year 2006

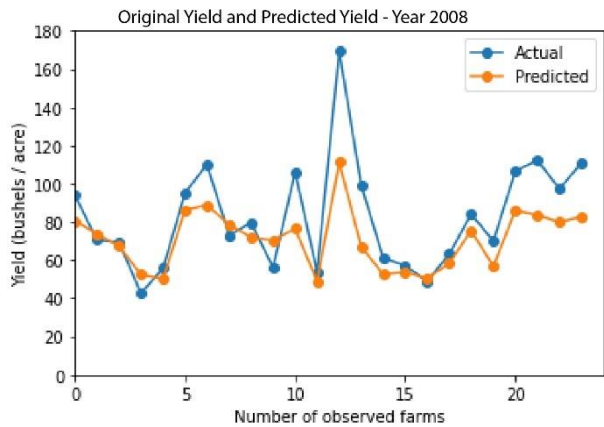




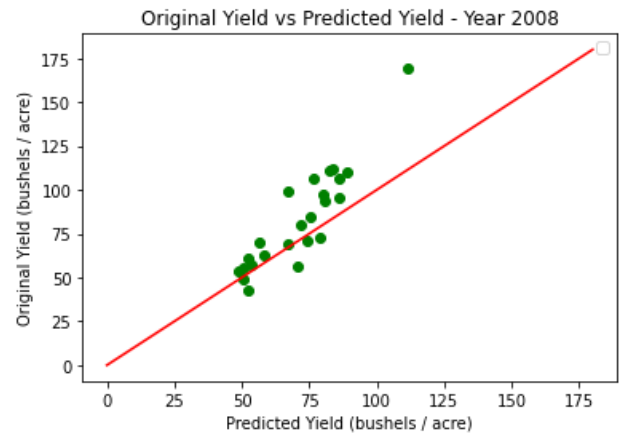
**Fig 4.20** Original yield and Predicted yield for the year 2007



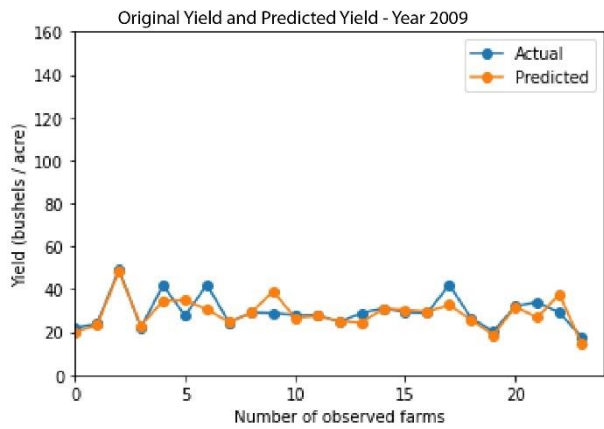
**Fig 4.21** Original yield vs Predicted yield for the year 2007



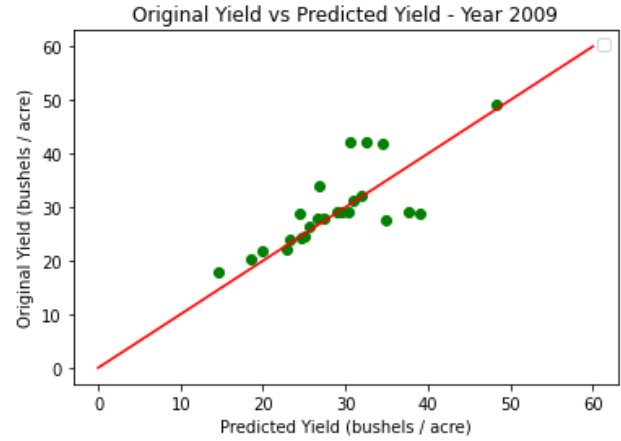
**Fig 4.22** Original yield and Predicted yield for the year 2008



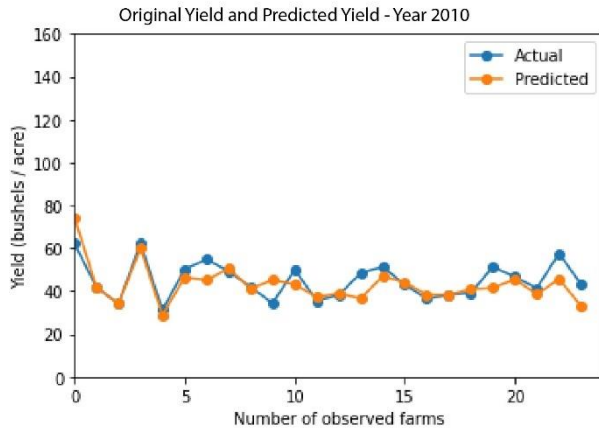
**Fig 4.23** Original yield vs Predicted yield for the year 2008



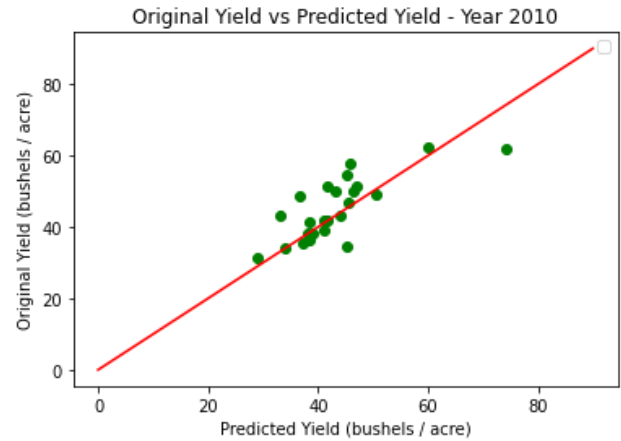
**Fig 4.24** Original yield and Predicted yield for the year 2009



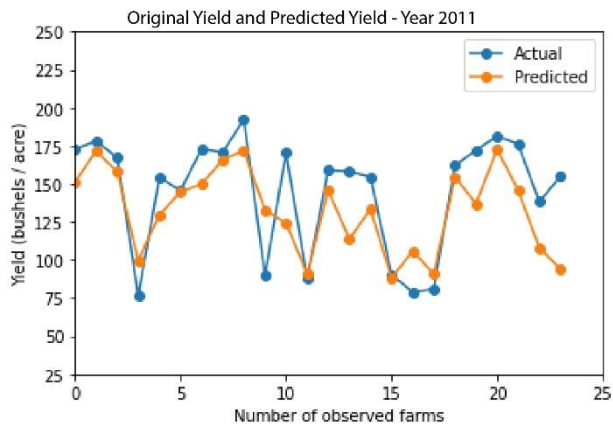
**Fig 4.25** Original yield vs Predicted yield for the year 2009



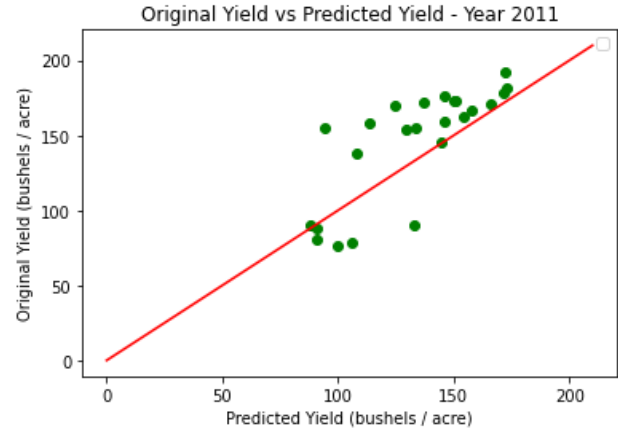
**Fig 4.26** Original yield and Predicted yield for the year 2010



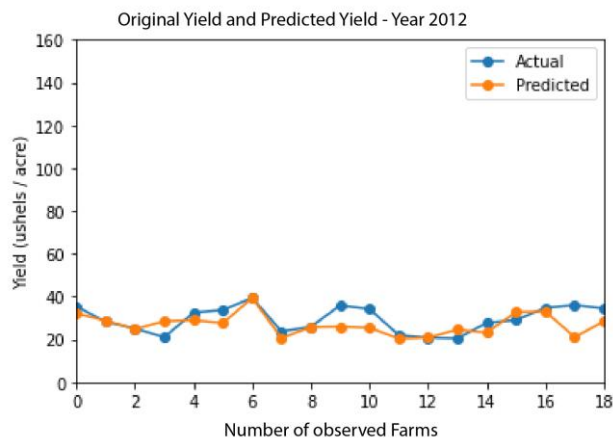
**Fig 4.27** Original yield vs Predicted yield for the year 2010



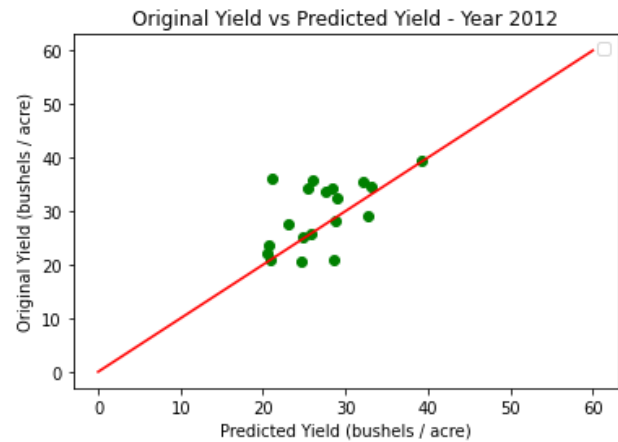
**Fig 4.28** Original yield and Predicted yield for the year 2011



**Fig 4.29** Original yield vs Predicted yield for the year 2011



**Fig 4.30** Original yield and Predicted yield for the year 2012



**Fig 4.31** Original yield vs Predicted yield for the year 2012

**Table 4.7** Year-wise analysis

Year	RMSE	MAPE	Full farm accuracy
2001	5.44	9.2	90.79
2003	5.8	18.54	81.5
2004	7.13	11.16	88.8

2005	18.05	14.22	85.77
2006	3.81	6.85	93.14
2007	7.74	9.97	90.02
2008	20.08	15.44	84.5
2009	4.42	10.42	89.58
2010	6.51	9.74	90.26
2011	24.7	14.32	85.6
2012	5.436	14.01	85.9

#### **4.3 DISCUSSIONS:**

1. Crop yield depends on various factors other than the ones that have been considered in the study, like soil quality, applied nitrogen quantity, plant diseases etc. These are factors that can't be represented by an aerial image. As observed in the year 2003, accuracy was the lowest compared to all the other years since it was a year “without summer” in Michigan which caused nitrogen deficiencies and weed problems in farms. Moreover soybean (crop grown in 2003) is a warm season crop, and these two factors had affected the yield of the crop. Nitrogen deficiencies and weed problems are features which weren't taken into account while training the model, hence the accuracy is about 81.5%. On the contrary, the year 2006 (crop grown Soybean) has a higher accuracy of around 93% since nitrogen deficiencies and weed problems weren't the important factors contributing to the yield.
2. Soil moisture profile / soil quality is an important parameter when it comes to yield prediction but this can't be represented with satellite images. Hence, it is observed that in the year 2005, where the crop grown is corn (a warm seasonal crop) and very dry , warm weather prevailed in Michigan state in mid-march. But due to the dry situation, the full soil moisture profile and the soil quality wasn't as expected and the crop yields were

affected due to the early start of the warm and dry conditions. Hence, the model had produced only an accuracy of 85.7% in the year 2005. Similarly, in the year 2008, Michigan experienced Hurricane Gustav in the month of September, leading to flooding and water-logged soils. Since the harvest of corn was in the month of October, there was a moderate impact on the structure of soil and yield produced. Hence, our model has produced an accuracy of about 84% in the year 2008.

3. As observed from table 4.7, wheat has a consistent accuracy among all the years 2001,2004,2007 and 2010 which is about 90%. It implies that the model has studied and can identify the pattern growth of wheat from plantation to harvest in a better manner compared to that of soybean and corn.
4. Similar results are shown in table 4.6, where crop-wise analysis was performed and wheat shows the highest accuracy of about 89.7%. This aligns with the conclusions stated above.

## **CHAPTER 5**

### **CONCLUSIONS AND FUTURE WORK**

In this work, two regression deep learning models were built to predict the crop yield in farm-level. The Landsat images of the study area, were cropped in a sliding four-pixel window to increase the size of the dataset. One of the models was trained using multi temporal spectral bands dataset while the other model using multi temporal vegetation indices dataset. The total farm's crop yield results of the spectral bands deep learning model provided a better accuracy of 94.3% compared to the accuracy provided by the vegetation indices which was 88.8%. Thus, the spectral band dataset was chosen as the features for the analysis of year-wise and crop-wise deep learning regression models. In this way our proposed model managed to obtain a good accuracy overcoming the challenges of predicting the crop yield for a small spatial farm area using remote satellite dataset.

The crop yield prediction's accuracy can be further increased even during the cases of weather anomaly by including the weather data parameters along with the spectral band features.

The crop yield dataset used in the proposed model, is in the duration between 2001-2011 (excluding 2002). These years don't have access to landsat-8 and landsat-9 satellite images which has a lesser cloud coverage compared to landsat-5 and landsat-7 satellite images that has been used in the proposed model. This implies, if the crop yield dataset for the farms is obtained for the years after 2013, there will be better access to more cloud free landsat-8 satellite images which in turn helps to create a model which can predict the crop yield data much before the harvest time.

## REFERENCES

- [1]Wolfert, Sjaak, Lan Ge, Cor Verdouw, and Marc-Jeroen Bogaardt. "Big data in smart farming—a review." *Agricultural systems* 153 (2017): 69-80.
- [2]Kamilaris, Andreas, Andreas Kartakoullis, and Francesc X. Prenafeta-Boldú. "A review on the practice of big data analysis in agriculture." *Computers and Electronics in Agriculture* 143 (2017): 23-37.
- [3]Schlenker, Wolfram, and Michael J. Roberts. "Nonlinear temperature effects indicate severe damages to US crop yields under climate change." *Proceedings of the National Academy of sciences* 106, no. 37 (2009): 15594-15598.
- [4]Makowski D. and Micheal L., “Use of dynamic model for predicting crop yield trends in foresight studies on food Security”, FAO Sixth International Conference On Agricultural Statistics ICAS VI, 2013
- [5]Taylor, J.C.; Wood, G.A.; Thomas, G. Mapping yield potential with remote sensing. *Precis. Agric.* 1997, 1, 713-720.
- [6] R Delrcolle, S J Maas, M Gurrit, and F Baret, “Remote sensing and crop production models : present trends,” *Journal of Photogrammetry and Remote Sensing*, vol. 47, no. 1986, pp. 145–161, 1992
- [7]Kim, Nari, and Yang-Won Lee. "Machine learning approaches to corn yield estimation using satellite images and climate data: a case of Iowa State." *Journal of the Korean Society of Surveying, Geodesy, Photogrammetry and Cartography* 34, no. 4 (2016): 383-390.
- [8]Terliksiz, Anil Suat, and D. Turgay Altýlar. "Use of deep neural networks for crop yield prediction: A case study of soybean yield in lauderdale county, alabama, usa." In *2019 8th international conference on Agro-Geoinformatics (Agro-Geoinformatics)*, pp. 1-4. IEEE, 2019
- [9]Arumugam, Ponraj, Abel Chemura, Bernhard Schauburger, and Christoph Gornott. "Remote Sensing Based Yield Estimation of Rice (*Oryza Sativa* L.) Using Gradient Boosted Regression in India." *Remote Sensing* 13, no. 12 (2021): 2379.
- [10]Aghighi, Hossein, Mohsen Azadbakht, Davoud Ashourloo, Hamid Salehi Shahrabi, and Soheil Radiom. "Machine learning regression techniques for the silage maize yield



- prediction using time-series images of Landsat 8 OLI." *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 11, no. 12 (2018): 4563-4577.
- [11]Kuwata, Kentaro, and Ryosuke Shibasaki. "Estimating crop yields with deep learning and remotely sensed data." In *2015 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, pp. 858-861. IEEE, 2015.
- [12]Ghazaryan, Gohar, Sergii Skakun, Simon König, Ehsan Eyshi Rezaei, Stefan Siebert, and Olena Dubovyk. "Crop yield estimation using multi-source satellite image series and deep learning." In *IGARSS 2020-2020 IEEE International Geoscience and Remote Sensing Symposium*, pp. 5163-5166. IEEE, 2020.
- [13]Engen, Martin, Erik Sandø, Benjamin Lucas Oscar Sjølander, Simon Arenberg, Rashmi Gupta, and Morten Goodwin. "Farm-Scale Crop Yield Prediction from Multi-Temporal Data Using Deep Hybrid Neural Networks." *Agronomy* 11, no. 12 (2021): 2576.
- [14]Russello, Helena. "Convolutional neural networks for crop yield prediction using satellite images." *IBM Center for Advanced Studies* (2018).
- [15]Jiang, Hao, Hao Hu, Renhai Zhong, Jinfan Xu, Jialu Xu, Jingfeng Huang, Shaowen Wang, Yibin Ying, and Tao Lin. "A deep learning approach to conflating heterogeneous geospatial data for corn yield estimation: A case study of the US Corn Belt at the county level." *Global change biology* 26, no. 3 (2020): 1754-1766.
- [16]Panda, Sudhanshu Sekhar, Daniel P. Ames, and Suranjan Panigrahi. "Application of vegetation indices for agricultural crop yield prediction using neural network techniques." *Remote Sensing* 2, no. 3 (2010): 673-696.
- [17] Battude, Marjorie, Ahmad Al Bitar, David Morin, Jérôme Cros, Mireille Huc, Claire Marais Sicre, Valérie Le Dantec, and Valérie Demarez. "Estimating maize biomass and yield over large areas using high spatial and temporal resolution Sentinel-2 like remote sensing data." *Remote Sensing of Environment* 184 (2016): 668-681.
- [18]Adeniyi, Odunayo David, Andrea Szabó, János Tamás, and Attila Nagy. "Wheat Yield Forecasting Based on Landsat NDVI and SAVI Time Series." (2020).
- [19]Magri, Antoni, Harold M. Van Es, Michael A. Glos, and William J. Cox. "Soil test, aerial image and yield data as inputs for site-specific fertility and hybrid management under maize." *Precision agriculture* 6, no. 1 (2005): 87-110.

[20]Thiam, S.; Eastmen R.J. Chapter on vegetation indices. In Guide to GIS and Image Processing, Volume 2; Idrisi Production: Clarke University, Worcester, MA, USA, 1999; pp. 107-122.

## **PUBLICATIONS**

1. Paper Title: Farm Level Crop Yield Prediction using Deep Learning with  
Multitemporal Satellite Images

Authors: Aruneswari S, Ramakrishna, Chandhana S, Nirali M Dave, Prathiksha R R,  
Dr. Aravinth J

Status: Drafting