

# Using Multilayer Perceptron (MLP) to predict crop yields

Thomas Donnelly  
Winona State University

tommymdonnelly@gmail.com

## ABSTRACT

This study aims to develop a Multilayer Perceptron (MLP) that accurately predicts crop yields within 10% of ground truth in 80% of cases using weather data, region, soil type, temperature, fertilizer, irrigation, days taken to harvest, and rail fall. The dataset has 1 million unique data points and 10 columns. The model will be trained using a random selection of 80% of the data for training and 10% for testing. The effectiveness of the model will be derived from the accuracy and the Mean Square Error (MSE) of the model.

## Categories and Subject Descriptors

D.3.3 [Programming Languages]: Language Constructs and Features – *abstract data types, polymorphism, control structures*. This is just an example, please use the correct category and subject descriptors for your submission.

## General Terms

Your general terms must be any of the following 16 designated terms: Algorithms, Management, Measurement, Documentation, Performance, Design, Economics, Reliability, Experimentation, Security, Human Factors, Standardization, Languages, Theory, Legal Aspects, Verification.

## Keywords

Keywords are your own designated keywords.

## 1. INTRODUCTION

As technology advances, it is important to explore as many ways it can be used to improve the world and the lives of people around us. One area where this could have a great impact is farming. Farming is notoriously difficult to predict. On average, crop yields in the US have increased since 1995, with 2023 yielding 179.3 bushels of corn per acre, compared to 113.5 in 1995 (USDA Crop, 1). This, however, can vary year to year depending on how much fertilizer is used, yearly rainfall, and the general region and climate. To help with this prediction, this study plans to create a neural network to assist in the prediction. This study plans to create a multilayer perceptron (MLP) to achieve this goal.

## 2. Hypothesis

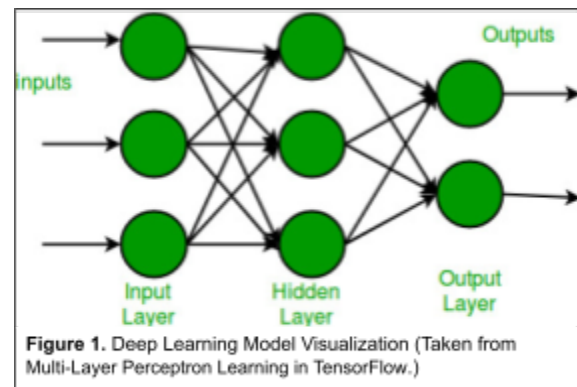
Multilayer Perceptron (MLP) accurately predicts crop yields within 10% of ground truth in 80% of cases using weather data, region, soil type, temperature, fertilizer, irrigation, days taken to harvest, and rail fall.

## 3. Methods

### 3.1 Model Basics

This project will focus on developing software. This study plans to use Python, a coding language with a large library of neural network creation tools, to create the model using the Pandas library. The model I will be creating is a Multilayer Perceptron (MLP). An MLP is made up of three layers. The first layer receives all input data, including weather data, region, etc. The second layer is all the fillings or the hidden layers. These

processes take the input data and try to make predictions based on it. While additional hidden layers can improve performance, having more does not necessarily result in a better model. Lastly, the output layer takes all the info the hidden layers give it and makes a final prediction or guess on the correct answer. See Figure 1 for reference.



The basic structure of the model using TensorFlow is as follows:

1. `Input(shape=(X_train.shape[1],))`. Input Layer that just ensures the model fits the data.
2. `Dense(x)`. `x` is the number of neurons that one layer can have to process data. This can be thought of as the thinking layer.
3. `BatchNormalization()`. This makes it so multiple Dense layers can read each other's output better.
4. `ReLU()` helps the model learn complex patterns in the data.
5. Repeat 2-4 with `x` either increasing or decreasing between layers.
6. `Dense(1)`. This layer is the final output and prediction.

### 3.2 Geometric Pyramid Formula

A formula was used to help narrow down the size of the model for training. This is called the Geometric Pyramid Formula. In short, the formula helps prevent creating a model that is too large for the data set given. The version of the formula used goes as follows:  $\# \text{ of parameters} / 10 > (\text{layer size} * \text{previous layer size}) + \dots + (\text{layer}(i) * \text{layer}(i - 1))$ .

### 3.3 Dataset

The dataset that will be used is taken from Kaggle, a site where you can get free datasets to work with. The dataset is titled "Agriculture Crop Yield Dataset" and is a 1 million non-null unique data point set (Agriculture Crop Yield Dataset). From my analysis of the data, it appears to be a generated set, which means it is not taken from the real world but instead created using similar data or trends in the field of interest. This is important to keep in mind as it may affect how the model can predict real-world data. The dataset has ten columns, which are as follows, with date types: Region object, Soil\_Type object, Crop object, Railfall\_mm float64, Temperature\_Celsius float64, Fertilizer\_Used boolean, Irrigation\_used boolean, Weather\_Condition object,

Days\_to\_Harvest int64, Yield\_tons\_per\_hectare float64 (Agriculture Crop Yield Dataset). The first nine columns will be used to predict the ground truth or actual value of the yield per hectare. Figures 2.1-2.5 are visualizations of the data.



Figure 2.1 Visualization of Regions in the Dataset (Created using Matplotlib)

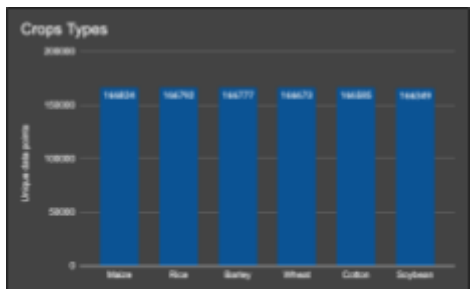


Figure 2.2 Number of Crop Types in the Dataset (Created using Matplotlib)

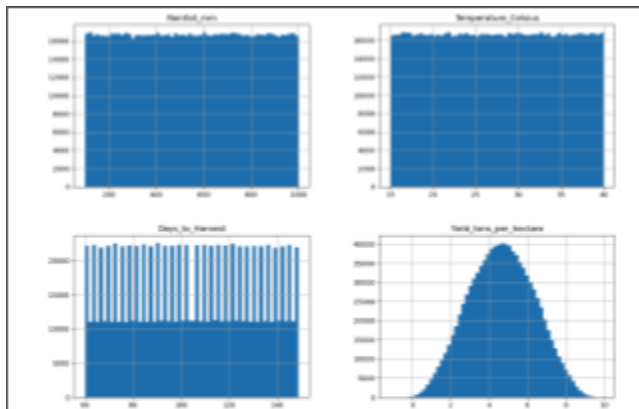


Figure 2.3 Display of the spread of Rainfall, Temperature, Days to Harvest, and Yield values. (Created using Matplotlib)

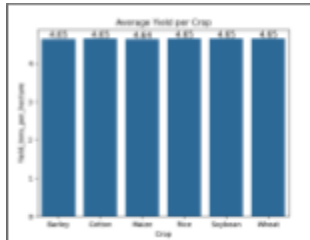


Figure 2.4 The average yield per crop type (Created using Matplotlib)

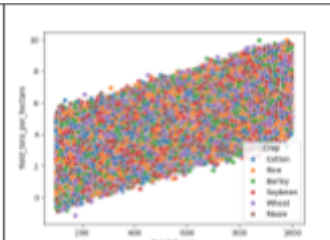


Figure 2.5 How rainfall affects crop yield (Created using Matplotlib)

### 3.4 Training

The method I will use for the training of the MLP is to divide the dataset into two parts. For the first part, 90% of the data will be available to the model to train on and evaluate its answers. Meanwhile, a random 10% will be set aside for after the model is done training. This data will be used to test how good the model is on previously unseen data. If the model was tested on data it had already seen, the chances would be high that it would recognize the data and learn the answer, and not how to get the answer. This is also assisted by having so many points of data, with 1 million being a large set, making it hard for the model to learn the answer and not the how.

### 3.5 Evaluation

Two methods will be used to show how effective the model is during and after it is created. The first is accuracy, the goal is for the model to at least 80% of the time be either 10% over or under the testing data. This can be implemented using the following logic::  $\text{accuracy} = \text{average}(\text{predicted} \geq \text{actual} * 0.9) \text{ OR } (\text{predicted} \leq \text{actual} * 1.1) \geq 0.8$ . The other method is Mean Square Error (MSE), this formula helps determine how accurate a model is over data, with a lower value meaning the more accurate it is (Mean Squared Error). See Figure 3 below for a similar research using MSE to display the accuracy of models. Their study is a good reference point for the project, as they tested models using only select columns to see which helped get to the ground truth the best. Meanwhile, this study will be using all the columns.

Table 5. Multiple factorial regression (NFR) models for different type of covariables.

Source <sup>1</sup>	df	Sum of squares ( $\times 10^3$ )	Mean squares ( $\times 10^3$ )	F	P
Year $\times$ treatment	207	279.52	13.503	5.69	0.0000
NET					
Treat $\times$ MEA	23	56.02	24.35	10.47	0.0000
Treat $\times$ MEF	23	54.05	23.45	9.89	0.0000
Treat $\times$ MEF2	23	46.26	20.12	8.79	0.0000
Treat $\times$ MEFM	23	13.39	5.82	2.49	0.0000
Deviation	115	99.09	8.616	3.97	0.0000
mT					
Treat $\times$ mTF	23	78.53	34.14	14.17	0.0000
Treat $\times$ mTF2	23	40.59	17.65	7.31	0.0000
Treat $\times$ mTF3	23	53.87	23.42	9.89	0.0000
Treat $\times$ mTFM	23	25.19	10.95	4.54	0.0000
Treat $\times$ mTFA	23	25.34	11.02	4.57	0.0000
Deviation	92	76.99	8.371	3.43	0.0000
mTU					
Treat $\times$ mTUF	23	71.81	31.22	12.85	0.0000
Treat $\times$ mTUF2	23	48.65	21.15	8.78	0.0000
Treat $\times$ mTUF3	23	27.52	11.96	4.96	0.0000
Treat $\times$ mTUFM	23	14.39	6.26	2.59	0.0000
Deviation	115	107.38	9.329	3.87	0.0000
PH					
Treat $\times$ PHD	23	16.88	7.34	2.94	0.0000
Treat $\times$ PHF	23	23.44	10.19	4.23	0.0000
Treat $\times$ PHM	23	27.18	11.82	4.89	0.0000
Treat $\times$ PHF2	23	14.49	6.30	2.61	0.0000
Deviation	115	182.68	15.88	6.38	0.0000
SH					
Treat $\times$ SHD	23	28.73	12.49	5.08	0.0000
Treat $\times$ SHF	23	15.33	6.66	2.70	0.0000
Treat $\times$ SHM	23	20.63	8.970	3.70	0.0000
Deviation	138	214.88	15.57	6.45	0.0000
EV					
Treat $\times$ EVD	23	69.68	30.30	12.81	0.0000
Treat $\times$ EVF	23	40.09	17.43	7.21	0.0000
Treat $\times$ EVM	23	19.19	8.34	3.53	0.0000
Treat $\times$ EVF2	23	18.31	7.979	3.31	0.0000
Treat $\times$ EVF3	23	14.08	6.126	2.54	0.0000
Deviation	92	138.88	15.18	6.31	0.0000
Error	409	118.879	2.906		

Figure 3. A similar study that compared multiple models and select data with Mean Square Error being one of the evaluations (Vargas 949)

## 4. EXPERIMENT

### 4.1 Data Normalization

The first problem encountered in creating a model was that a third of the data could not be run due to being in an unsupported

format. Figure 4 below shows the columns divided by their data formats.

Categorical	Numerical	Boolean
Region	Rainfall_mm	Fertilizer_Used
Soil_Type	Temperature_Celsius	Irrigation_Used
Crop	Days_to_Harvest	
Weather_Condition		

Figure 4. Training columns sorted by data type

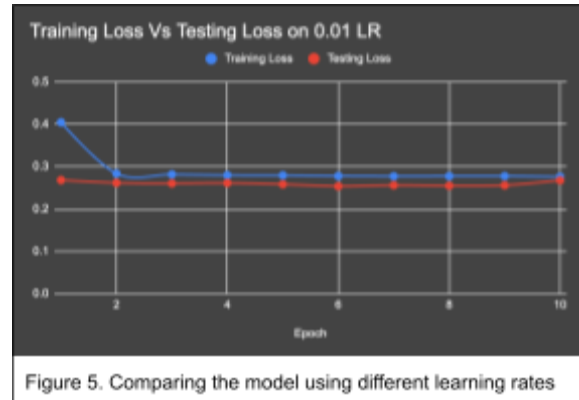
Sklearn was used to create a pipeline to separate and manage the data in the format it's in. Since Region, Soil\_Type, Crop, and Weather\_Condition were strings, they could be run as is. The boolean or yes-no columns Fertilizer\_Used and Irrigation\_Used had to be transformed to be processed. The method used to achieve this was to convert the data to numbers, with True becoming 1 while False became 0. All the numerical values were able to be run, but a normalizer was applied to help the model read the data, and if missing values were found, the average score of the set was put in their place. After all the data was modified, it was all put back together, and the model was ready to start running.

## 4.2 Hardware

The expectation was to run the model on a dedicated NVIDIA GPU for its training. This would allow for quicker training, a larger possible model size, and a larger batch size. A larger model size would allow for more neurons; this does not always mean more performance out of a model but would allow the ability to see if a bigger model would be better. The batch size is how much data in a run-over data is stored in memory at a time. A larger batch size generally means higher performance of a model. It was found, however, that the versions of Nvidia GPU driver (Nvidia-smi) and Compute Unified Device Architecture (NVCC) that were installed on the machine were too new for Tensorflow, a library used to create the model, and a different machine was chosen. In the end, a MacBook M2 was used to train the model. This meant that the model had to be trained on a CPU, but given the time constraints, it was good enough.

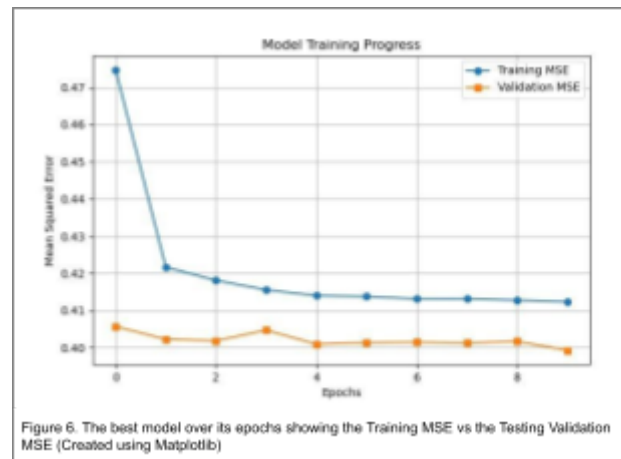
## 4.3 Learning Rate

The learning rate, or the rate at which the model can learn the was chosen to be 0.001. A rate of 0.01 was tested but was not used due to the model learning the training data and learning how to predict using the data too quickly. Figure 5 is a quick example of running the model on a 0.01 learning rate. As the epochs increase or the number of times the model goes over the training data, the rate of loss on the testing data starts to reverse. Loss is a measure of how confident the model is in its answer. The reverse in loss in the testing data, but not a similar reverse in training data, shows that the model started to learn the data and not from the data.



## 5. RESULTS

See Figure 6 for the best run as it was trained over its 10 Epochs with a 0.001 learning rate and a best of 0.3979 Mean Squared Error. Figure 7 shows how the best result is structured. It has 10 Epochs, 4 hidden layers with their size in order being 256, 128, 64, and 32 neurons.



## 6. ANALYSIS

Using the best model, the study was able to achieve a total of 62.46% of the predictions being within 10% of the ground truth yield and 87.98 % of the predictions being within 20% of the ground truth yield. This has not met the goal of 80% within 10% of the ground. Training the model with the best parameters over longer epochs would not yield better results, as the model would start to learn the set and could lose accuracy. If you see Figure 8, you can see how little the MAE decreases as the epochs increase. The parameters do, however, have room to increase due to the model being well under the budget of 100,000. The model is 45,825 using the Geometric Pyramid Formula. Though when tested, it was found to slightly decrease the accuracy.



Epoch	Training Loss	Training MAE	Validation Loss	Validation MAE
1	0.9848	0.6191	0.2589	0.4057
2	0.2812	0.423	0.2542	0.4023
3	0.2754	0.4185	0.2538	0.4019
4	0.272	0.4161	0.2575	0.4047
5	0.2693	0.4138	0.2528	0.401
6	0.269	0.4137	0.2535	0.4014
7	0.2682	0.4129	0.2536	0.4015
8	0.2693	0.414	0.2533	0.4013
9	0.2682	0.4132	0.2536	0.4017
10	0.268	0.4129	0.2507	0.3993
Final Evaluation			0.2485	0.3978

Figure 8. Table showing the best model's training data

## 7. CONCLUSION

In conclusion, the study did not achieve the goal it set out to achieve. If the goal post was moved to just 20% within the ground truth, then it would have achieved such. Meanwhile, 62.26% is not nothing. I believe this study has achieved its goal using a network to predict crop yield given the data.

## 8. REFERENCES

- [1] Kaggle. "Agriculture Crop Yield Dataset." Kaggle, 2021, <https://www.kaggle.com/datasets/samuelotiattakorah/agriculture-crop-yield/data>.
- [2] Jha, Dinesh, et al. "Application of Artificial Neural Networks for Predicting Crop Yield: A Review." *Frontiers in Plant Science*, vol. 10, 2019, <https://www.frontiersin.org/journals/plant-science/articles/10.3389/fpls.2019.00621/full>.
- [3] "Corn Yield." USDA National Agricultural Statistics Service, [https://www.nass.usda.gov/Charts\\_and\\_Maps/graphics/cornyield.pdf](https://www.nass.usda.gov/Charts_and_Maps/graphics/cornyield.pdf).
- [4] "Crop Progress & Condition." USDA National Agricultural Statistics Service, 2024, [https://www.nass.usda.gov/Charts\\_and\\_Maps/Crop\\_Progress\\_and\\_Condition/2024/index.php](https://www.nass.usda.gov/Charts_and_Maps/Crop_Progress_and_Condition/2024/index.php).
- [5] Subramanian, S., et al. "Early Crop Yield Prediction for Agricultural Decision Making Using Remote Sensing Data." *Journal of Water and Climate Change*, vol. 14, no. 12, 2023, pp. 4729-4741, <https://iwaponline.com/jwcc/article/14/12/4729/99202/Early-crop-yield-prediction-for-agricultural>.
- [6] Vargas, Mateo. "Interpreting Treatment × Environment Interaction in Agronomy Trials." *Agronomy Journal*, vol. 93, no. 4, 2001, pp. 949-960.
- [7] "Multi-Layer Perceptron Learning in TensorFlow." GeeksforGeeks, <https://www.geeksforgeeks.org/multi-layer-perceptron-learning-in-tensorflow/>.
- [8] "Mean Squared Error." GeeksforGeeks, <https://www.geeksforgeeks.org/mean-squared-error/>.
- [9]