# WEEKLY PROGRESS REPORT

**Name:** Sahil Sanjay Pokharkar

**Domain:** Data Science and Machine Learning

**Date of Submission:** 28/01/2026

**Project Title:** Prediction of Crop Production in India

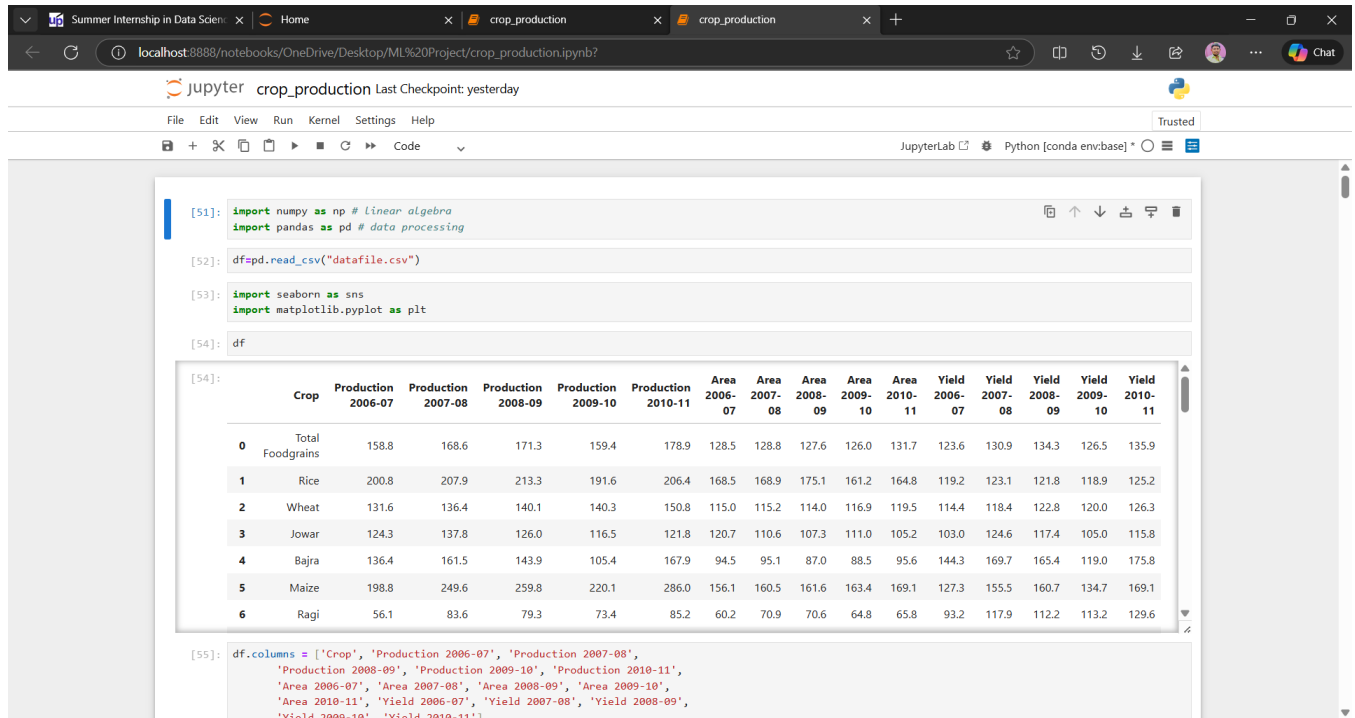**Organization:** Upskill Campus

**Week Ending:** Week 03

# ACKNOWLEDGEMENT

# Table Of Content

# Chapter 1: Overview

The third week of the internship focused on **feature engineering, data encoding, and implementation of initial machine learning models** for predicting crop production. The objective was to transform the pre-processed dataset into a machine-learning-ready format and build baseline models to evaluate prediction performance.

# Chapter 2: Tasks Completed & Achievements

**1. Feature Engineering and Encoding**

- Performed encoding of categorical variables such as State, District, Crop, and Season using suitable encoding techniques.
- Normalized and scaled numerical features like Area to improve model performance.
- Finalized the feature set required for training machine learning models.

```
df.columns = ['Crop', 'Production 2006-07', 'Production 2007-08',
        'Production 2008-09', 'Production 2009-10', 'Production 2010-11',
        'Area 2006-07', 'Area 2007-08', 'Area 2008-09', 'Area 2009-10',
        'Area 2010-11', 'Yield 2006-07', 'Yield 2007-08', 'Yield 2008-09',
        'Yield 2009-10', 'Yield 2010-11']
```

**2. Dataset Splitting**

- Split the dataset into training and testing sets to evaluate model generalization.
- Ensured appropriate data distribution between training and testing data.

```
df.columns = df.columns.str.strip()
```

```
df.shape
```

```
(55, 16)
```

```
df.duplicated().sum()
```

```
np.int64(0)
```

**3. Machine Learning Model Implementation**

- Implemented baseline machine learning models such as:
- Linear Regression
- Decision Tree Regressor
- Trained models on the processed dataset to predict crop production values.

```
X = df.drop(["Crop", "Yield 2010-11"], axis=1)
y = df["Yield 2010-11"]
```

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42
)
```

```
from sklearn.linear_model import LinearRegression

model = LinearRegression()
model.fit(X_train, y_train)
```

```
▼ LinearRegression  ⓘ ?

LinearRegression()
```

## 4. Model Evaluation

- Evaluated model performance using metrics such as:
- Mean Absolute Error (MAE)
- Mean Squared Error (MSE)
- R² Score
- Compared model results to identify the better-performing algorithm.

```
Mean Squared Error: 117.29278407595895
R2 Score: 0.8800067327354562
```

## 5. Performance Analysis

- Analysed overfitting and underfitting issues in initial models.
- Observed how feature selection and scaling impacted prediction accuracy.



Actual vs Predicted Yield

# Chapter 3: Challenges

- **High Cardinality in Categorical Features:** Encoding state and district columns increased feature dimensionality.
- **Overfitting in Decision Tree Model:** Required careful tuning of model parameters.
- **Computational Complexity:** Training models on large datasets increased processing time.
- **Metric Interpretation:** Understanding regression evaluation metrics required additional practice.

```
[7]: df.shape
```

```
[7]: (55, 16)
```

# Chapter 4: Lessons Learned

- Gained hands-on experience in converting real-world data into machine-learning-ready format.

- Learned how different algorithms behave on the same dataset.

- Developed a strong understanding of regression evaluation metrics.

- Understood the importance of feature scaling and encoding in predictive modelling.

- Improved debugging and optimization skills during model training.



Feature Correlation Heatmap

# Chapter 5: Learning Resources

- Scikit-learn official documentation
- Online tutorials on regression algorithms
- Upskill Campus learning materials and mentor guidance
- Blogs and articles on crop yield prediction using machine learning

# Chapter 6: Next Week's Goals

- Implement advanced models such as **Random Forest and Gradient Boosting**.
- Perform hyperparameter tuning to improve model performance.
- Apply cross-validation techniques.
- Compare all models and select the best performing one.
- Prepare the model for final deployment and reporting.

# Chapter 7: Additional Comments

Week 03 was a crucial phase in transitioning from data analysis to machine learning implementation. Building and evaluating predictive models provided valuable practical exposure to real-world data science workflows. This week significantly strengthened my confidence in applying machine learning techniques to solve agricultural prediction problems.