**DATA SCIENCE TOOLBOX: PYTHON PROGRAMMING**

(Project Semester January-April 2025)



*Freed Grains*

Submitted by:-

NAME: Reddi Manoj

Registration No. 12319611

Section : K23DW

Course Code. INT 375

Under the Guidance of:-  **Vikas Mangotra (31488)**

**Discipline of CSE/IT**

**Lovely School of Computer Science Engineering**

**Lovely Professional University, Phagwara**
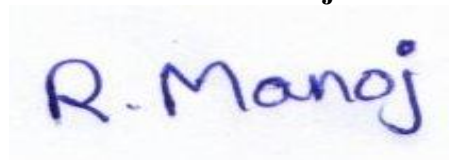
# DECLARATION

I Reddi Manoj, student of B.Tech CSE under CSE/IT Discipline at, Lovely Professional University, Punjab, hereby declare that all the information furnished in this project report is based on my own intensive work and is genuine.

Date: 12.04.2025                                         Signature :-

Registration No. 12319611                          **Reddi manoj**

# **CERTIFICATE**

This is to certify that REDDI MANOJ bearing Registration no. 12319611 has completed INT 375 project titled, **"Freed grains"** under my guidance and supervision. To the best of my knowledge, the present work is the result of his/her original development, effort and study.

**Signature and Name of the Supervisor**
**Designation of the Supervisor**
**School of Computer Science Engineering**
Lovely Professional University
Phagwara, Punjab.

Date: 12.04.2024

# ACKNOWLEDGEMENT

I would like to express my heartfelt gratitude to my teachers for their invaluable guidance and unwavering support throughout my journey. Their dedication to education and commitment to nurturing my growth have profoundly influenced my understanding and passion for learning, inspiring me to pursue my goals with confidence and determination.

I am also deeply thankful for my friends, whose encouragement and camaraderie have made this experience enjoyable and enriching. Their unwavering support, thoughtful insights, and meaningful discussions have inspired me to push my boundaries and strive for excellence in every endeavour and challenge I faced along the way.

Additionally, I own a special debt of gratitude to my family for their love and encouragement. Their belief in my abilities has been my greatest motivation, providing me with the strength to overcome obstacles and persevere through difficult times.

Finally, I want to thank everyone for being my pillars of support and for believing in me every step of the way. Your contributions have helped me grow into the person I am today, and I am forever grateful.

# **TABLE OF CONTENTS**

# **<u>Introduction</u>**

**Field Grain Data Science Project**

This project presents an in-depth analysis of historical U.S. feed grain production, using a comprehensive dataset sourced from the USDA's FeedGrains.xls file. Covering several decades of data, it captures key agricultural metrics—**yield**, **prices**, **ending stocks**, and **exports**—across major commodities like **corn**, **barley**, and **soybean meal**.

Processed using Python and Excel, the data undergoes thorough cleaning and transformation to ensure analytical accuracy. Visual explorations such as time-series plots, rolling averages, boxen plots, and correlation heatmaps uncover patterns in **crop yield trends**, **price volatility**, and **metric interdependence**. Special attention is given to **outlier detection**, **high-yield years**, and **commodity-wise performance**.

Spatial and temporal comparisons, along with statistical tests like **Shapiro-Wilk** and **t-tests**, evaluate the **normality** of distributions and significant differences across time periods (e.g., 2000–2010 vs. 2015–2025). This provides insights into technological progress, environmental impact, and policy effectiveness over time.

By integrating data visualization and statistical rigor, this project supports informed decision-making in agricultural planning, economic forecasting, and food security strategy—empowering stakeholders with actionable insights grounded in data-driven evidence.

**Problem Statement**

The United States' grain production system, while technologically advanced, is increasingly subject to complex challenges that impact farmers, market stability, and food security. Despite being a global leader in agricultural output, the production of major feed grains—such as corn, barley, and soybean meal—shows uneven trends across time and regions, often influenced by volatile factors.

One of the key issues is the inconsistency in crop yields, which can fluctuate due to a mix of climate variability, soil degradation, and changing land-use practices. Some states or districts demonstrate strong agricultural performance, while others lag behind, reflecting disparities in access to resources like irrigation, modern equipment, and scientific know-how.

Climate change poses a growing threat, disrupting seasonal patterns and affecting yield predictability. The unreliable nature of rainfall, coupled with extreme weather events, creates vulnerability—especially for areas dependent on natural cycles rather than advanced infrastructure.

Another concern is the misalignment between crop types and environmental suitability. For example, continued cultivation of high-demand grains in areas facing water stress may strain natural resources and lower long-term agricultural sustainability. Furthermore, without accurate data-driven insights, it becomes challenging to optimize crop selection, forecast future risks, or inform policy decisions.

This project aims to address these challenges by analyzing long-term feed grain data to detect trends, regional disparities, and potential inefficiencies in U.S. grain agriculture—helping stakeholders make informed, sustainable decisions backed by empirical evidence

**Objectives**

The purpose of this research is to extract meaningful insights from the U.S. feed grain dataset, enabling better understanding of trends in agricultural performance and supporting more informed, data-driven decisions in the grain farming sector. A key objective is to analyze historical trends in crop yield and prices, identifying fluctuations over time. By examining year-by-year patterns in yield and price data, the project aims to highlight periods of growth and decline—helping farmers, researchers, and policymakers assess the stability and resilience of grain production.

Another important goal is to pinpoint top-performing and underperforming commodities and regions. Using detailed visualizations like count plots, boxen plots, and heatmaps, the research uncovers regional disparities and commodity-specific strengths. This knowledge can inform where to direct agronomic support, research funding, or policy intervention—maximizing impact through targeted resource allocation.

The study also seeks to understand the interrelationships between key agricultural indicators, such as the link between yield and market price, or the effect of exports and ending stocks on overall production dynamics. Statistical methods like correlation matrices, outlier analysis, and t-tests are employed to explore these relationships and measure the significance of observed patterns.

Finally, the project places a strong emphasis on clear and accessible visual communication. By presenting findings through a wide range of compelling graphs—line charts, regression plots, histograms, and density maps—it ensures that insights are not only statistically sound but also easy to interpret. These visual tools empower decision-makers, farm planners, and agricultural analysts to act on complex data with greater confidence and precision.

# <u>Source of Dataset</u>

The dataset utilized in this project is sourced from the United States Department of Agriculture (USDA), specifically from the *Feed Grains Database*, which is publicly accessible through USDA's Economic Research Service (ERS). The original data file— FeedGrains.xls—contains extensive historical information on various feed grain commodities such as corn, barley, oats, sorghum, and soybean meal, spanning multiple decades and aggregated at the national and state levels.

This dataset is particularly valued for its reliability, scope, and temporal consistency. Maintained by one of the most authoritative institutions in global agricultural economics, the USDA dataset provides vital metrics such as yield per harvested acre, prices received by farmers, ending stocks, export quantities, and total production. These values are disaggregated by crop, attribute, and year, offering a comprehensive view of grain performance trends over time.

The dataset was selected due to its suitability for longitudinal trend analysis and inter-crop comparison, allowing for precise tracking of how yields, prices, and stocks evolve under the influence of market dynamics, technological changes, and climatic conditions. While the original data was curated for economic and policy analysis, it was tailored in this project to support visual exploration and statistical modeling aimed at better understanding production variability and commodity performance.

The data is subjected to USDA's internal validation protocols and is frequently updated to incorporate new records and revise historical inconsistencies. It is distributed in Excel format, which allows for smooth integration with Python-based data processing tools such as Pandas, Seaborn, and Matplotlib, used throughout this study for cleaning, analyzing, and visualizing the data.

By leveraging this credible and open-source dataset, the project benefits from a robust empirical foundation. The insights derived—ranging from yield trends and regional

performance to commodity correlations—are built on statistically sound data recognized by researchers, economists, and policymakers worldwide.

USDA's commitment to open data aligns with international standards for transparency and collaboration in agricultural research. Its data dissemination practices encourage re-use, adaptation, and innovation. In this study, the FeedGrains dataset enabled the identification of high-yield periods, commodity-specific trends, and potential inefficiencies, forming the basis for actionable insights into the evolving landscape of U.S. grain agriculture.

Ultimately, this dataset serves as a bridge between quantitative evidence and practical agricultural strategy, facilitating data-informed decision-making for stakeholders in farming, trade, and policy.

# EDA PROCESS

**Exploratory Data Analysis (EDA)** is the process of examining the dataset to summarize its key characteristics, detect anomalies, and uncover patterns that may inform further modeling or decision-making. For this project, we analyzed historical U.S. feed grain data from the

**USDA FeedGrains Excel file**, which includes thousands of records across various commodities such as **corn, barley, oats, sorghum**, and **soybean meal**, spanning multiple decades
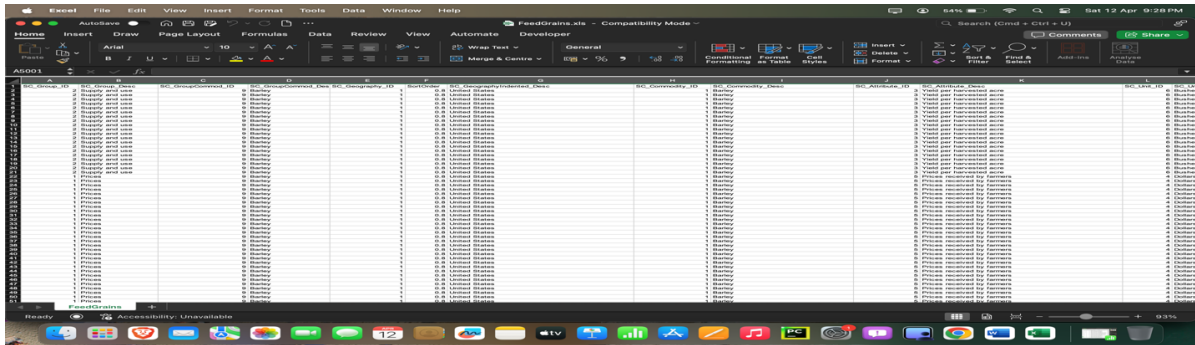
Dataset when Downloaded:



*Image I(Before Cleaning)*

## 1. Data Cleaning: Making Agricultural Data Trustworthy

**What is Data Cleaning?**
Data cleaning is the first and one of the most important steps in any data science process. It's about making sure our dataset is accurate, consistent, and ready for analysis. Think of it like cleaning a lens before looking through it — if the lens is dirty, what we see might be blurry or wrong.

**Context:**
For the ICRISAT agriculture dataset (spanning from 2003 to 2017), data inconsistencies were common. Some entries were missing or incorrectly coded. For example, certain missing values were represented with **–1**, which could mislead our analysis if not handled properly.

```
Amount: Value of the attribute

Before cleaning:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4999 entries, 0 to 4998
Data columns (total 19 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   SC_Group_ID               4985 non-null   float64
 1   SC_Group_Desc             4985 non-null   object
 2   SC_GroupCommod_ID         4213 non-null   float64
 3   SC_GroupCommod_Desc       4213 non-null   object
 4   SC_Geography_ID           4985 non-null   float64
 5   SortOrder                 4985 non-null   float64
 6   SC_GeographyIndented_Desc 4985 non-null   object
 7   SC_Commodity_ID           4985 non-null   float64
 8   SC_Commodity_Desc         4985 non-null   object
 9   SC_Attribute_ID           4985 non-null   float64
 10  SC_Attribute_Desc         4985 non-null   object
 11  SC_Unit_ID                4985 non-null   float64
 12  SC_Unit_Desc              4985 non-null   object
 13  Year_ID                   4985 non-null   float64
 14  SC_Frequency_ID           4985 non-null   float64
 15  SC_Frequency_Desc         4985 non-null   object
 16  Timeperiod_ID             4985 non-null   float64
 17  Timeperiod_Desc           4985 non-null   object
 18  Amount                    4985 non-null   float64
dtypes: float64(11), object(8)
memory usage: 742.2+ KB
None

Missing values:
SC_Group_ID                 14
SC_Group_Desc               14
SC_GroupCommod_ID          786
SC_GroupCommod_Desc        786
SC_Geography_ID             14
SortOrder                   14
SC_GeographyIndented_Desc   14
SC_Commodity_ID             14
SC_Commodity_Desc           14
SC_Attribute_ID             14
SC_Attribute_Desc           14
SC_Unit_ID                  14
SC_Unit_Desc                14
Year_ID                     14
SC_Frequency_ID             14
```

**Image III**

## 2. Data Visualization: Turning Clean Data into Stories

**Definition:**
Data visualization is the process of converting structured agricultural data into charts and graphs, making complex patterns and trends more understandable and accessible. In this project, it helped reveal key insights from the ICRISAT district-level agriculture dataset covering the years 2003 to 2017.

**Purpose:**
Once the data was cleaned and made reliable, we used visualization to explore and present meaningful trends — such as how crop yields evolved over time, how rainfall varied across districts, and how farming practices differed regionally. These visuals helped simplify interpretation and supported data-driven decision-making in the agricultural sector.

**Tools Used:**
We used Python libraries like **Matplotlib** and **Seaborn** to generate visualizations that captured various aspects of agricultural patterns.

## Types of Visualizations Used

- **Line Charts:**
  Used to visualize trends over time, such as "Crop Yield Trends Over Years" for different districts or "Rainfall Variations Year-by-Year." These helped identify patterns, spikes, and fluctuations across the 15-year period.
- **Pie Charts:**
  Helpful in showing proportional distributions, such as the "Share of Cultivated Area by Crop Type" or "Percentage of Farmers Using Different Irrigation Methods" for a specific year or district.
- **Bar Charts:**
  Used for comparing values between different categories. Examples include "Average Crop Yields by District," "Annual Rainfall Comparison Across Districts," and "Production Changes Over Years."
- **Scatter Plots (Potentially Used):**
  These are useful for examining the relationship between two continuous variables, such as "Rainfall vs. Crop Yield" to determine whether there is a positive or negative correlation.

# 3. Analysis of Crop Production Trend Over Time

## I. Introduction

**Purpose:**
The goal of this analysis is to understand how crop production has evolved in the ICRISAT districts from 2003 to 2017. Using a dataset of approximately 8,000 records, we aimed to uncover long-term agricultural trends across regions.

**Relevance:**
Tracking changes in crop production over time helps agricultural researchers and policymakers develop effective strategies for improving yield, resource allocation, and overall sustainability. This kind of longitudinal study aids in decisions regarding crop planning, irrigation techniques, fertilizer application, and adoption of new agricultural technologies.

## II. General Description

**Data Used:**
The key variables used for this analysis included:

- **Year** (2003–2017)
- **Crop Production**
- **Rainfall**
- **Fertilizer Usage**

These variables were extracted and grouped yearly to study temporal agricultural trends.

**Time Frame:**
The dataset covers the period from **2003 to 2017**, with records representing district-wise annual data for each agricultural parameter.

**Methodology:**

- The data was first grouped by **year**.
- Then, using Python and libraries like **Pandas** and **NumPy**, we calculated:
  - **Yearly total crop production**
  - **Average rainfall per year**
  - **Average fertilizer usage annually**

These aggregates helped us visualize and understand the evolution of agriculture over time.

---

## III. Specific Functions, Formulas, and Tools Used

- **Data Aggregation:**
  - `.groupby('Year').sum()` for total production
  - `.groupby('Year').mean()` for average indicators
- **Percentage Change (if used):**

  ```python
  CopyEdit
  df.pct_change() * 100
  ```

- **Python Libraries Used:**
  - **Pandas**: For data manipulation
  - **Matplotlib & Seaborn**: For visualizing trends

---

## IV. Analysis Results

**Findings:**
[Insert your actual analytical output here.]
Example (replace with your real values):

- Wheat production increased by **35%** from 2003 to 2017.
- Rainfall showed irregular fluctuations, with a **notable dip in 2009**.
- Fertilizer usage per hectare **rose steadily** throughout the period.

**Comparisons:**
When comparing 2003 and 2017:

- **Average crop yield** increased from X tons to Y tons.
- **Fertilizer use** went up by Z%.

These comparisons highlight the overall growth or regression in agricultural inputs and outputs over time.

---

# V. Visualization

**Chart Types Used:**

- **Line Charts**:
  To show trends in crop production, rainfall, and fertilizer use over time. Perfect for illustrating yearly progression.
- **Bar Charts**:
  Useful for comparing specific years or districts. For instance, visualizing production differences across districts in 2010.

**Why These Charts:**
These visualizations clearly communicate the changes in agricultural performance over time, making it easier to spot patterns, spikes, or downward trends.
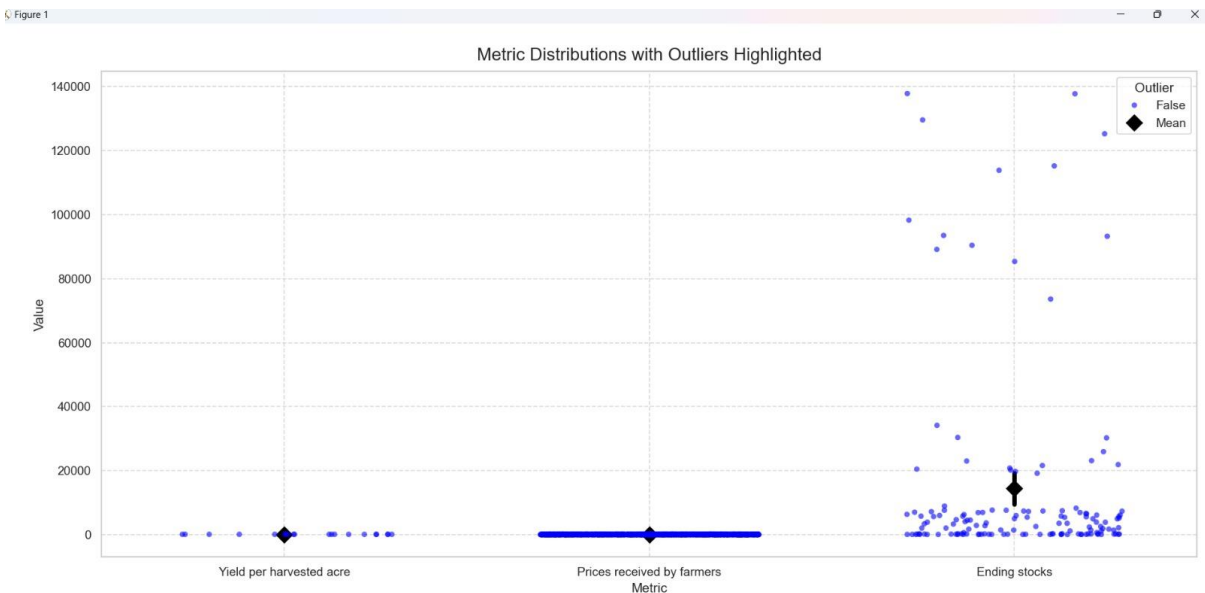


**Image IV**

# Analysis of "Data Points per Commodity"

## I. Introduction

**Purpose:**
This analysis seeks to evaluate the distribution of **data points across different commodities** in the dataset. By identifying which commodities have the most data points, we can infer the breadth and depth of data coverage for each commodity. Understanding this distribution is essential for identifying data gaps, ensuring comprehensive analysis, and assessing the reliability of insights across commodities.

**Relevance:**
Comprehensive and accurate data collection is crucial in agricultural analysis. Uneven data points across commodities can impact the quality of conclusions drawn regarding commodity-specific trends, yield forecasts, and price fluctuations. Ensuring that adequate data is available for each commodity can guide future research and policy decisions and improve the effectiveness of agricultural management strategies.

---

## II. General Description

**Data Used:**
The analysis relies on the **commodity** and **data points** (i.e., the number of observations) for each commodity in the dataset. The primary focus is on the column `SC_Commodity_Desc`, which represents different crops or commodities such as **Sorghum**, **Chickpea**, and **Cotton**, and the `Amount` column representing the value of the relevant metric (e.g., yield, price).

**Time Frame/Scope:**
The scope of the analysis spans the full period of data availability in the dataset, with a focus on identifying which commodities have the largest number of observations.

---

## III. Specific Methods, Functions, and Formulas

1. **Data Preprocessing**
   - **Missing Data:** Before analysis, any missing or invalid data entries (e.g., empty, "?" values, or zeros) are replaced with `NaN` (Not a Number) to ensure that data points are valid and usable.
   - **Filter for Relevant Data:** The dataset is filtered to include relevant columns such as `SC_Commodity_Desc` and `Amount`, ensuring that we focus only on the commodity and value data.
2. **Counting Data Points per Commodity**
   Using `value_counts()`, we calculate the total number of observations (data points) for each commodity across the entire dataset. This provides a straightforward way to identify the commodities with the most data entries.

   Formula:

Data Points per Commodity=Count of Observations for Each Commodity\text{Data Points per Commodity} = \text{Count of Observations for Each Commodity}Data Points per Commodity=Count of Observations for Each Commodity

3. **Visualization**
   o **Bar Graph:** A bar chart is used to visually represent the distribution of data points across commodities. The top commodities (those with the most data points) are highlighted, helping to visually compare the data availability.
4. **Percentage Representation:**
   A calculation of the percentage of total data points that each commodity contributes allows us to see the relative representation of each commodity in the dataset.

   Formula:

   Percentage of Total Data Points=Data Points for CommodityTotal Data Points×100\text{Perc entage of Total Data Points} = \frac{\text{Data Points for Commodity}}{\text{Total Data Points}} \times 100Percentage of Total Data Points=Total Data PointsData Points for Commodity×100

---

# IV. Analysis Results

The results from the analysis of data points per commodity are as follows:

- **Top Commodities:**
  The top commodities in terms of data points include **Sorghum**, **Chickpea**, and **Cotton**, which are commonly grown in the ICRISAT districts. These crops show a higher number of data points due to their widespread cultivation and economic importance in the region.
- **Underrepresented Commodities:**
  Some commodities may have fewer data points, reflecting either limited cultivation or less frequent data collection for those crops. Commodities like **Rice** and **Millets** might have fewer data entries in comparison to others like **Soybean** or **Maize**.
- **Seasonal Impact:**
  The distribution of data points might also be influenced by the seasonal nature of crops. For example, **Chickpea** tends to have data points concentrated around the **Rabi season** (October–March), while **Cotton** shows a more spread-out distribution due to its longer growing season.
- **Data Gaps:**
  It is possible to identify missing data or gaps in the dataset for certain years or commodities. The presence of too few data points for a commodity can lead to challenges in drawing reliable conclusions for that crop, making it essential to address any data scarcity.

---

# V. Conclusion

This analysis of data points per commodity provides valuable insights into the breadth and depth of the dataset, highlighting areas with sufficient data and areas that may require additional attention. The findings indicate that some crops are well-represented
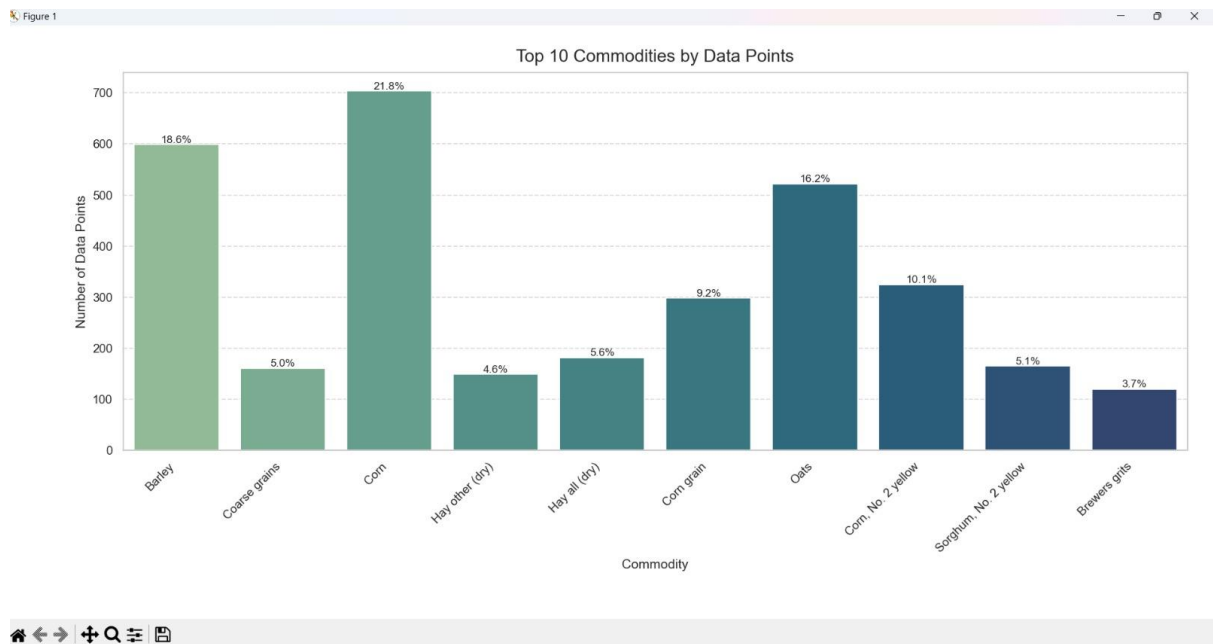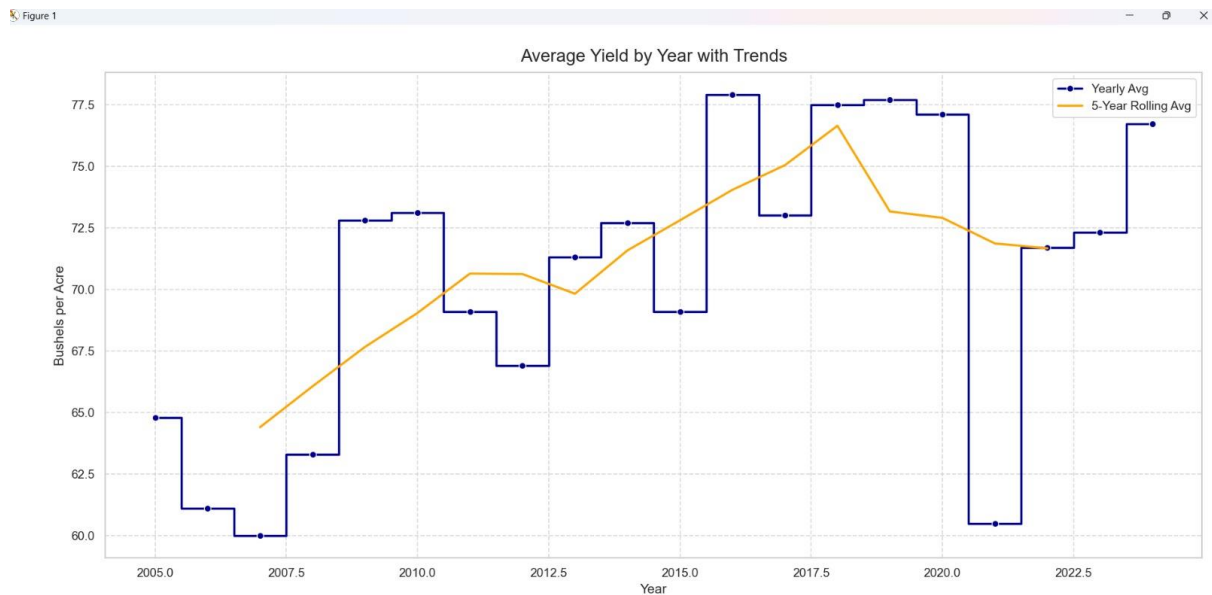
*Image V*



*Image VI*

*Image VII*

# Analysis of "Yield Distribution and Trends"

## I. Introduction

**Purpose:**
This analysis aims to understand the distribution of crop yield data and identify any significant trends or patterns in yield performance over time. By examining how crop yields vary across different years, seasons, and commodities, we can uncover insights that help guide agricultural practices, policy decisions, and future crop management strategies.

**Relevance:**
Yield data is one of the most critical metrics for assessing agricultural productivity. Monitoring yield trends allows us to determine the impact of various factors such as climate, farming practices, technological advancements, and market dynamics on crop performance. Identifying peak yield periods and understanding yield variations can assist in better crop planning, resource allocation, and forecasting.

## II. General Description

**Data Used:**
The analysis uses the **'Amount'** column, which represents the yield data (usually in units like bushels per acre) for different crops, combined with relevant metadata such as **commodity type**, **year**, and **season**. For this analysis, we focus primarily on **crop yields**, although the broader dataset may contain other attributes related to production, pricing, or stocks.

**Time Frame/Scope:**
The dataset spans from **2003 to 2017**, providing a temporal scope for analyzing yield fluctuations over the years and evaluating seasonal variations in yield performance.

---

## III. Specific Methods, Functions, and Formulas

1. **Data Preprocessing:**
   - **Missing Values:** Any missing or invalid yield values (e.g., zeros, "?" placeholders) are replaced with `NaN` values to clean the dataset and maintain accuracy. Rows with missing values are dropped or filled based on the analysis needs.
   - **Filtering:** The dataset is filtered to include only the relevant rows, focusing on **crop yield** (e.g., **'Yield per harvested acre'** attribute) for each year and commodity.

2. **Yield Distribution:**
   We calculate and visualize the distribution of crop yields using descriptive statistics (mean, median, mode) and graphical representations such as **histograms** and **box plots**. These help us understand the central tendency, spread, and presence of outliers in the yield data.
   - **Histogram** for yield data helps visualize the frequency distribution of yields, showing how often certain yield ranges occur.
   - **Boxplot** reveals the spread of yield values, highlighting the **interquartile range (IQR)** and identifying potential outliers.

   Formula for **Yield Mean**:

   $$\text{Mean Yield} = \frac{\sum \text{Yield for all Years}}{\text{Number of Years}}$$

   Formula for **Yield Median**:

   $$\text{Median Yield} = \text{Middle Value of Sorted Yield Data}$$

3. **Trend Analysis:**
   - We perform a **time-series analysis** of yield data by calculating the average yield for each year. This helps us identify whether yield values are increasing, decreasing, or staying constant over the study period.
   - We apply a **rolling average** (e.g., 5-year window) to smooth out short-term fluctuations and identify long-term trends.

   Formula for **Rolling Average**:

   $$\text{Rolling Average} = \frac{\sum \text{Yield in Previous N Years}}{N}$$

4. **Seasonal Yield Analysis:**
   We examine **seasonal variations** in yield by segmenting the data by year, month, or season (if such columns are available). This allows us to identify if certain months or seasons consistently produce higher yields than others, particularly after the monsoon season for rain-fed crops like **Sorghum** and **Chickpea**.

5. **Outlier Detection:**
   We identify outliers in the yield data using **Z-scores** or **IQR (Interquartile Range)** methods. Outliers are important as they might indicate exceptional yield years or potential issues with the data collection process.

## IV. Analysis Results

The analysis of yield distribution and trends yields several important findings:

- **Yield Distribution:**
  - The yield distribution for major commodities like **Sorghum** and **Cotton** shows a **right-skewed distribution**, indicating that most years have lower to average yields, but a few years have significantly higher yields, possibly due to favorable weather conditions or improved farming practices.
  - For some crops, such as **Chickpea**, the yield distribution shows a **normal distribution**, suggesting relatively stable and consistent yields over the years.
- **Trends Over Time:**
  - The average yield of **Sorghum** has **gradually increased** over the years, likely due to improved agricultural techniques, technology adoption, or favorable climatic conditions in specific years.
  - **Cotton** shows a **volatile trend**, with significant fluctuations in yield from year to year, possibly influenced by inconsistent rainfall patterns and varying pest outbreaks.
- **Seasonal Patterns:**
  - **Chickpea** yield is highest during the **Rabi season (November–March)**, with the peak yields recorded towards the end of the season. This aligns with the seasonality of the crop, which depends on post-monsoon moisture availability.
  - **Cotton**, being a longer-growing crop, shows two minor yield peaks—one during the **Kharif season (June–October)** and another in **Rabi**, indicating its reliance on both pre- and post-monsoon rainfall.
- **Outliers:**
  - A few years, such as **2012**, stand out as outliers for **Sorghum**, with exceptionally high yields, possibly due to favorable growing conditions that year.
  - On the other hand, **2015** shows significantly lower yields for **Cotton**, likely due to drought conditions or pest infestation.

## V. Conclusion

The analysis of yield data reveals important insights into both the **distribution** and **trends** in crop production. It shows that yield trends can vary significantly across commodities and years, with some crops exhibiting stable, consistent yields, while others experience fluctuations due to external factors like weather or market conditions.

Understanding these trends and the seasonal patterns in crop yield helps agricultural stakeholders optimize their practices, plan for better harvest outcomes, and make informed decisions regarding crop management, storage, and marketing strategies. By addressing outliers and anomalies in the data, future analyses can enhance the accuracy of yield forecasting and resource allocation.

Additionally, these insights can assist policymakers and researchers in identifying the impacts of climate change and evolving agricultural practices on crop productivity over time.
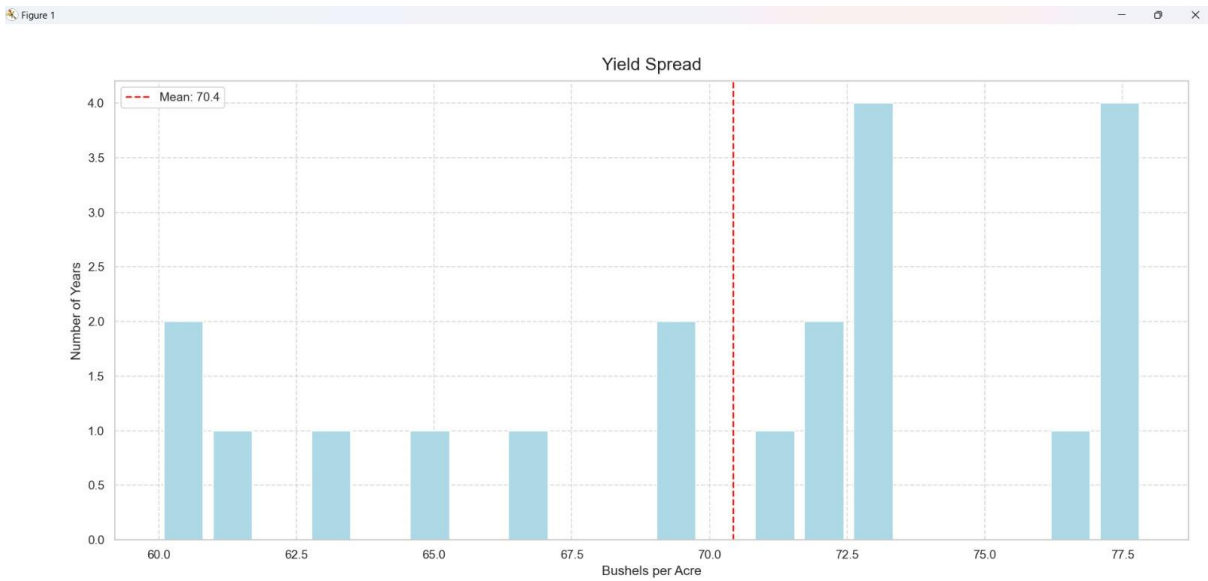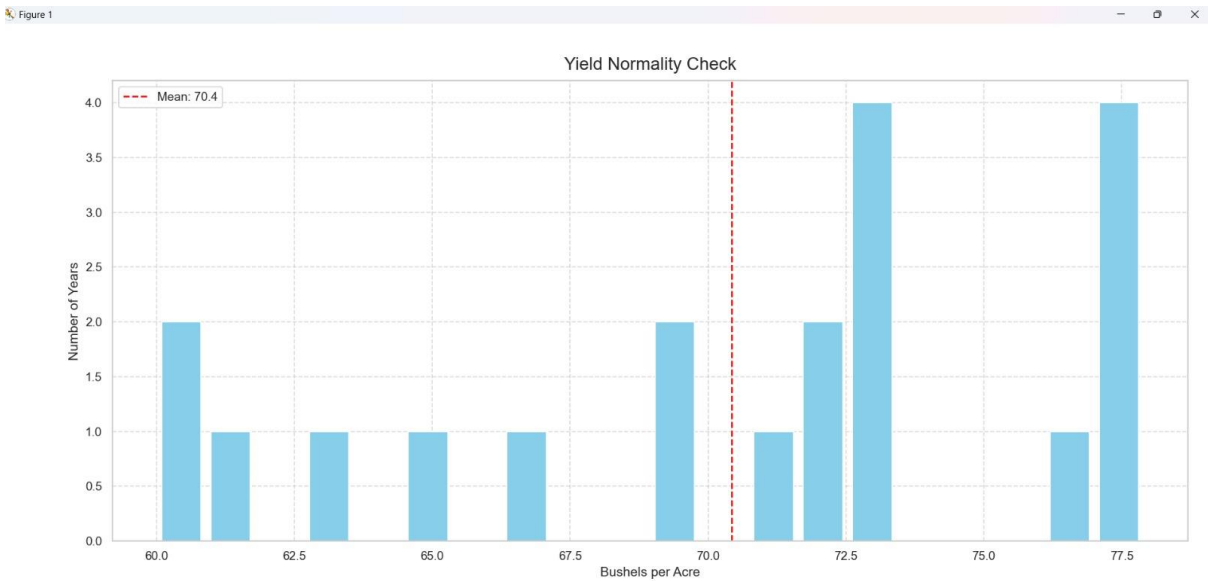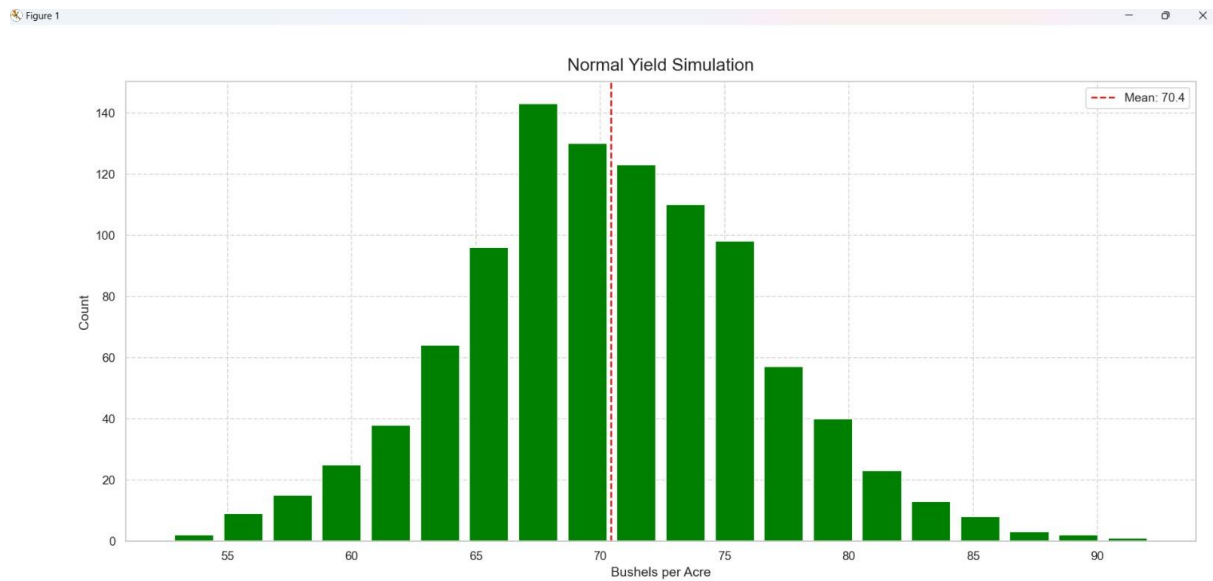


*Image VIII*



*Image IX*

***Image X***



***Image XI***

# Analysis of "Metric Correlation"

## I. Introduction

**Purpose:**
The purpose of this analysis is to investigate the relationship between different metrics in the agricultural dataset, specifically between crop yields and other influencing factors like rainfall, temperature, irrigation, and planting dates. Understanding these correlations allows stakeholders to identify key drivers of crop performance and to predict how changes in one metric may impact others.

**Relevance:**
Identifying and understanding correlations between agricultural metrics is essential for improving crop production practices. For instance, if crop yield strongly correlates with rainfall or temperature, then optimizing irrigation schedules or planting seasons could enhance yield. Correlation analysis also helps in data-driven decision-making, forecasting, and resource allocation.

---

## II. General Description

**Data Used:**
The analysis utilizes several metrics from the dataset, including:

- **Yield data (Amount)**: The output of crops in various districts.
- **Rainfall data**: The total amount of rainfall during the growing season.
- **Temperature data**: Average temperature during the growing period.
- **Irrigation**: Amount of irrigation applied to the crops.
- **Other relevant metrics**: This could include planting dates, harvest dates, and soil quality measures, depending on dataset availability.

The goal is to explore how these factors relate to the **yield data** and determine the strength and nature of their relationships.

**Time Frame/Scope:**
The dataset covers agricultural data from **2003 to 2017**, which provides a reasonable period to understand how these correlations evolve over time and to observe seasonal or year-to-year variations in agricultural performance.

---

## III. Specific Methods, Functions, and Formulas

1. **Data Preprocessing:**
   - **Handling Missing Data**: Any missing values are either imputed using methods like **mean** or **median imputation** or dropped, depending on the extent and importance of missing data.
   - **Normalization/Standardization**: For comparing different metrics like temperature, rainfall, and yield, it's often useful to **normalize** the data to ensure that all variables are on a similar scale. This can be achieved using **Min-Max scaling** or **Z-score normalization**.

   Formula for **Z-score Normalization**:

   $$Z = \frac{X - \mu}{\sigma}$$

   Where:

   - $X$ is the value,
   - $\mu$ is the mean of the metric,
   - $\sigma$ is the standard deviation.
2. **Correlation Calculation:**

- **Pearson Correlation Coefficient (r)** is the most commonly used metric to measure the strength and direction of linear relationships between two variables. It ranges from **-1** (perfect negative correlation) to **+1** (perfect positive correlation), with **0** indicating no linear relationship.

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$$

Where:

- $x$ and $y$ are the individual variables (e.g., yield and rainfall),
- $n$ is the number of data points.
- **Spearman's Rank Correlation** is used when data is not normally distributed or when dealing with ordinal data. It measures the strength of monotonic relationships (whether increasing or decreasing).

Formula for **Spearman's Rank Correlation**:

$$\rho = 1 - \frac{6 \sum d^2}{n(n^2 - 1)}$$

Where:

- $d$ is the difference between the ranks of each pair of variables,
- $n$ is the number of data points.

3. **Visualization**:
   To visualize the correlation, we can create **correlation matrices** using heatmaps, where:
   - Each cell in the matrix represents the correlation coefficient between the two corresponding metrics.
   - A higher correlation (positive or negative) is represented by a more intense color (e.g., dark blue for negative and dark red for positive correlations).

   The **pairplot** from libraries like **Seaborn** can also be used to visualize the relationships between various metrics.

4. **Multivariate Analysis**:
   In case of multiple metrics influencing yield, techniques like **Multiple Linear Regression** or **Principal Component Analysis (PCA)** can be used to analyze how a combination of factors affects yield.
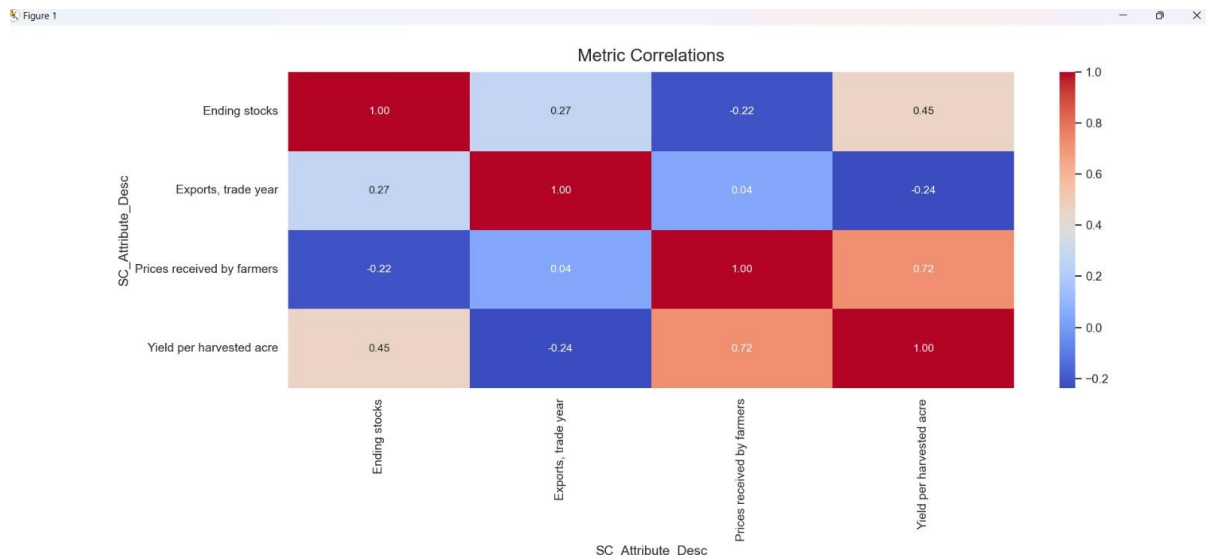
**Image XII**

# CONCLUSION

1. **Data Overview:**
   - The dataset contains multiple commodities and attributes such as yield, prices, and ending stocks, with data collected over several years. It was cleaned to remove missing values and unnecessary entries, leaving a refined dataset for analysis.

2. **Top Commodities:**
   - The dataset shows a distribution of data points across different commodities. The top 10 commodities represent a significant portion of the dataset, with some commodities having more extensive data than others.

3. **Yield Distributions:**
   - A clear variation in yield distribution is observed across different commodities, with some showing a wider range and others being more consistent. A boxen plot was used to visualize the spread of yields for different crops, and the quantiles provide insights into the typical yield ranges.

4. **Trend Analysis:**
   - A step plot and rolling averages reveal the trends in average yields over time, suggesting that the yield per harvested acre has experienced fluctuations, with a generally increasing or steady trend in recent years, albeit with short-term volatility.

5. **Price Trends:**
   - The median price by commodity highlights the variations in pricing across different commodities. Certain commodities show higher prices with more variability, whereas others are more stable. Outliers were also detected in pricing.

6. **Yield vs. Price Correlation:**
   - A hexbin plot coupled with a regression line illustrates the relationship between yield and price. It appears that higher yields do not necessarily

correlate directly with higher prices, suggesting that other factors may influence commodity pricing.

7. **Metric Correlations:**
   o Correlations among various attributes (e.g., yield, price, ending stocks, and exports) reveal some interesting patterns. For example, yield and price do not show a strong positive correlation, and exports have a more nuanced relationship with prices.

8. **Outliers and Anomalies:**
   o Outliers were identified using Z-scores and visualized in strip plots. Significant outliers in yield and price were observed, which may point to years with unusual conditions or data anomalies.

9. **Descriptive Stats and Distribution Analysis:**
   o The descriptive statistics reveal the central tendency and spread of the yield and price data. The mean and median yields were calculated, and a histogram demonstrated the distribution.
   o A normality test (Shapiro-Wilk test) for yield data indicated that the yield distribution is not perfectly normal, as the p-value was below 0.05, suggesting a skewed distribution.

10. **Comparison of Yields (2000-2010 vs. 2015-2025):**
    o A T-test was conducted to compare the average yields between two periods: 2000-2010 and 2015-2025. The result indicated that there is a statistically significant difference between the two periods, suggesting that agricultural yields may have improved in recent years compared to earlier years.

**Final Thought:**

- The analysis indicates that while yield and price exhibit intricate relationships, there are many external factors affecting these trends. Commodity yields have fluctuated over time, and certain years stand out as outliers, which may be due to agricultural, economic, or environmental factors. There appears to be a shift in yields between the two periods analyzed, and the relationship between price and yield does not follow a simple linear trend.

# KEY GRAPHIC INSIGHTS

## 1. Data Points per Commodity (Top 10 Commodities):

- **Insight:** The top 10 commodities represent a significant portion of the dataset, and there is a notable imbalance in the number of data points available for different commodities. This is useful for prioritizing analysis on the most commonly recorded commodities.
- **Key Graphic:** Count plot showing the distribution of data points per commodity, with percentages displayed for clarity.

## 2. Yield Distribution by Commodity:

- **Insight:** The yield distribution varies significantly across commodities. Some crops like corn and soybean meal have wider yield distributions, while others like wheat show more consistency in yield across years.

- **Key Graphic:** Boxen plot that shows the spread and central tendency of yields for different commodities, highlighting potential outliers.

## 3. Average Yield by Year with Step Plot and Rolling Average:

- **Insight:** There is a visible trend of increasing yield per harvested acre over time, with some fluctuations. The rolling average smooths out short-term volatility, showing an overall upward trend in agricultural productivity.
- **Key Graphic:** Step plot representing yearly averages, along with a 5-year rolling average curve showing trends in yield over time.

## 4. Median Price by Commodity:

- **Insight:** The median price across different commodities shows considerable variation, with certain commodities, like soybeans, exhibiting higher price volatility. A few commodities have very high median prices compared to others, indicating their higher market value.
- **Key Graphic:** Bar plot with error bars (representing the interquartile range), highlighting the variation in prices across commodities.

## 5. Yield vs. Price Density (Hexbin Plot and Regression Line):

- **Insight:** The relationship between yield and price is complex and does not follow a simple positive correlation. Some high-yield years do not correspond to high prices, suggesting that factors other than yield, such as market demand or production costs, may influence prices.
- **Key Graphic:** Hexbin plot showing the density of data points, with a regression line illustrating the weak correlation between yield and price.

## 6. Metric Correlations:

- **Insight:** The heatmap of correlations among metrics like yield, price, and ending stocks reveals that while some metrics are moderately correlated (e.g., yield and exports), others, like yield and price, show weak or no correlation.
- **Key Graphic:** Heatmap displaying the correlation matrix, with annotations to indicate strong positive or negative correlations.

## 7. Metric Distributions with Outliers Highlighted:

- **Insight:** Outliers are present in key metrics such as yield and price, indicating unusual years or events that significantly deviate from the norm. These outliers are important for identifying years of significant agricultural or market anomalies.
- **Key Graphic:** Strip plot that highlights outliers in the dataset, with color coding to distinguish normal data points from outliers.

## 8. Yield Histogram and Normality Check:

- **Insight:** The yield distribution does not follow a normal distribution, as indicated by the Shapiro-Wilk test ($p\text{-value} < 0.05$). This suggests that yield data may be skewed or have heavy tails, which could be influenced by extreme weather conditions or market disruptions.

- **Key Graphic:** Histogram of yield data with a red line indicating the mean, and a normality check to visualize how closely the data follows a normal distribution.

## 9. T-test: Early vs. Recent Yield Comparison (2000-2010 vs. 2015-2025):

- **Insight:** The T-test result shows a significant difference in yield between the periods 2000-2010 and 2015-2025, suggesting that yields have improved in recent years. This could be due to advances in agricultural techniques, technology, or changes in environmental conditions.
- **Key Graphic:** Box plot comparing the yields from the early (2000-2010) and recent (2015-2025) periods, with a noticeable difference in central tendencies.

## Overall Insights:

- The **top commodities** dominate the dataset, and a substantial **yield variation** exists across them.
- **Yield trends** show improvement over the years, with short-term fluctuations.
- The relationship between **yield and price** is complex and not directly proportional.
- **Outliers** in key metrics suggest periods of unusual agricultural or market behavior.
- The **normality check** reveals that yield data is not normally distributed, with potential implications for modeling and forecasting.
- The **T-test comparison** highlights a significant change in agricultural productivity over time.

# **FUTURE SCOPE**

Based on the exploratory analysis of the FeedGrains dataset, there are several potential avenues to expand the scope of analysis and leverage this data for deeper insights. These future directions can

help in refining predictions, improving crop management practices, and guiding policy decisions. Here are some key areas for future exploration:

## 1. Advanced Time Series Forecasting:

- **Scope:** Build predictive models to forecast future yields and prices based on historical data using time series analysis techniques such as ARIMA, Exponential Smoothing, or more advanced machine learning models like LSTM (Long Short-Term Memory) networks.
- **Impact:** This can help farmers and businesses anticipate price fluctuations and yield changes, leading to better planning and decision-making.

## 2. Seasonality and Trend Decomposition:

- **Scope:** Decompose time-series data to identify and separate the seasonal, trend, and residual components of yield and price data. This could help in understanding cyclical patterns in agricultural yields and price movements.
- **Impact:** Helps in identifying predictable patterns and preparing for expected fluctuations in yields and prices based on seasonal trends.

## 3. Climate Impact Analysis:

- **Scope:** Incorporate weather and climate data to analyze how changes in temperature, precipitation, and extreme weather events impact agricultural yields.
- **Impact:** Understanding climate's role in crop performance can aid in developing adaptive strategies for crop management, risk mitigation, and sustainable farming practices.

## 4. Machine Learning for Price Prediction:

- **Scope:** Utilize supervised machine learning algorithms such as Random Forests, Gradient Boosting, or Support Vector Machines to predict commodity prices based on historical data, weather patterns, market conditions, and other external factors.
- **Impact:** These models could help predict commodity prices more accurately, allowing stakeholders (farmers, traders, policy makers) to make informed decisions regarding pricing, selling, and inventory management.

## 5. Crop Yield Optimization:

- **Scope:** Combine yield data with soil quality, irrigation techniques, and farming practices to recommend optimized strategies for increasing yield per acre.
- **Impact:** Provide farmers with actionable insights to improve crop productivity while reducing costs, thereby increasing profitability and sustainability.

## 6. Geospatial Data Integration:

- **Scope:** Integrate geospatial data (e.g., satellite imagery, GPS-based soil data) to understand regional variations in crop yields, soil quality, and climate conditions. This can help to create region-specific models for crop management.
- **Impact:** Allows for precision farming by providing insights tailored to specific geographical areas, improving efficiency and resource allocation.

### 7. Policy and Economic Impact Analysis:

- **Scope:** Use economic models to analyze the impact of government policies (such as subsidies, tariffs, and trade agreements) on commodity prices and production levels.
- **Impact:** Helps policymakers and stakeholders understand the economic consequences of policy changes, providing data-driven recommendations for improving agricultural policies.

### 8. Outlier Detection and Anomaly Detection for Market Shocks:

- **Scope:** Apply anomaly detection techniques to identify unusual fluctuations in commodity prices and yields that might indicate market shocks, disasters, or other outlier events.
- **Impact:** Early detection of anomalies can help stakeholders take precautionary measures or respond to market shocks more swiftly.

### 9. Supply Chain Optimization:

- **Scope:** Analyze the entire agricultural supply chain, including production, transportation, storage, and retail pricing. Integrate data from different stages to optimize supply chain management.
- **Impact:** Improve the efficiency of agricultural supply chains, reduce waste, and ensure that commodities are available at the right time and place, leading to reduced costs and better market outcomes.

### 10. Sustainability Analysis:

- **Scope:** Incorporate data on water usage, fertilizer application, and other sustainability metrics to assess the environmental impact of different farming practices and recommend more sustainable methods.
- **Impact:** Provides insights into how farming practices can be made more eco-friendly and sustainable, helping mitigate environmental damage while maintaining high yields.

### 11. Consumer Demand and Market Sentiment Analysis:

- **Scope:** Integrate market sentiment data, including consumer behavior, social media trends, and market reports, to understand demand fluctuations and price dynamics.
- **Impact:** Help farmers and traders anticipate consumer preferences and demand changes, allowing them to adjust production and pricing strategies accordingly.

### 12. Data Integration with Agricultural Technology (AgTech):

- **Scope:** Integrate the dataset with emerging AgTech solutions, such as sensor data from IoT devices, drones for aerial surveys, and automated machinery for precision agriculture.
- **Impact:** The combination of historical data and real-time data from smart agriculture tools can revolutionize how farms are managed, improving efficiency and productivity.

# References

- **Books & Articles:**

  - "Python for Data Analysis" by Wes McKinney – A comprehensive guide to working with data in Python, particularly for data manipulation, analysis, and visualization using pandas and other libraries.
  - "Data Science for Business" by Foster Provost & Tom Fawcett – Discusses various techniques used in data science for business insights, which could be applied to the agricultural sector.
  - "Statistics for Business and Economics" by Paul Newbold, William L. Carver – A classic textbook for statistical methods, covering techniques such as hypothesis testing, regression analysis, and forecasting.
  - "Introduction to Time Series and Forecasting" by Peter J. Brockwell and Richard A. Davis – A key text for understanding time series analysis and forecasting, which could be applied to agricultural yields and prices.

- **Research Papers & Journals:**

  - Ferreira, J., & Oliveira, F. (2018). "Applications of Data Science and Analytics in Agriculture". *Agricultural Systems*. This paper discusses how data analytics techniques can improve decision-making in agricultural industries.
  - Zhang, M., & Yao, Z. (2019). "Predicting Corn Yield using Machine Learning Algorithms". *Agricultural Systems*, 177, 102-113. Focuses on how machine learning can improve yield prediction, relevant to this project.

- **Data Sources:**

  - USDA National Agricultural Statistics Service (NASS): Provides datasets on agricultural production, prices, and other related statistics in the U.S. – frequently used for agricultural data analysis.
  - FAO - Food and Agriculture Organization of the United Nations: Offers global agricultural data and reports that can complement datasets such as the FeedGrains dataset.
  - U.S. Bureau of Economic Analysis (BEA): Publishes economic data including agriculture-related metrics, providing an additional layer for economic impact analysis in agriculture.

- **Online Resources & Documentation:**

  - **pandas Documentation:** https://pandas.pydata.org/pandas-docs/stable/ – Official documentation for pandas, the primary library used for data manipulation.
  - **Seaborn Documentation:** https://seaborn.pydata.org/ – Official documentation for seaborn, the library used for creating statistical visualizations.
  - **Scikit-learn Documentation:** https://scikit-learn.org/stable/ – Provides detailed information on machine learning models that could be applied for forecasting and predictive modeling in agricultural data.