

DATA SCIENCE TOOLBOX: PYTHON PROGRAMMING (INT 375)

Sr. No.	Registration NO.	Name of Students	Roll No	Total Marks	Marks Obtained	Signature
1	12314228	SHRIJAN	66			



L OVELY
P ROFESSIONAL
U NIVERSITY

Transforming Education Transforming India

Submitted to Ms. Gargi Sharma

Lovely Professional University

Jalandhar, Punjab, India

Delivered by:	Received by:
Name of Students: Shrijan Reg no: 12314228,	Name of Faculty: Ms. Gargi Sharma UID: 29439 Signature:

DATA SCIENCE TOOLBOX: PYTHON PROGRAMMING (INT 375)

PROJECT REPORT

(Project Semester January-April 2025)

(Exploratory Data Analysis on Agricultural Land Usage in India)

Submitted by

(Shrijan)

Registration No12314228

Programme and Section:- CSE K23GR

Course Code INT375

Under the Guidance of

(Ms. Gargi Sharma)

Discipline of CSE/IT

Lovely School of Computer Science and Engineering

Lovely Professional University, Phagwara

CERTIFICATE

This is to certify that Shrijan bearing Registration no. 12314228 has completed INT375 project titled, “Exploratory Data Analysis on Agricultural Land Usage in India” 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 and Engineering

Lovely Professional University

Phagwara, Punjab.

Date: 12th April, 2025

DECLARATION

I, Shrijan, student of Computer Science and Engineering 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: 12th April 2025

Registration No. 12314228

Signature

Shrijan

A handwritten signature in black ink, appearing to read 'Shrijan', enclosed within a large, loopy oval stroke.

Acknowledgement

I would like to express my sincere gratitude to all those who contributed to the successful completion of this project.

First and foremost, I am deeply thankful to **Gargi Sharma**, my mentor and guide, for their constant support, valuable suggestions, and encouragement throughout the duration of this project. Their insights helped me explore the dataset more meaningfully and present my findings effectively.

I would also like to thank the Department of Computer Science and Engineering, **LOVELY PROFESSIONAL UNIVERSITY**, for providing me with the necessary resources and a learning environment that made this project possible.

My heartfelt thanks to my friends and classmates who supported me during the various stages of this analysis, offering valuable feedback and motivation.

Lastly, I am grateful to the creators and maintainers of the dataset, as well as the open-source community for tools like Python, Pandas, Seaborn, and Matplotlib, which played an integral role in my data analysis journey.

This project has been a great learning experience, and I truly appreciate everyone who helped make it possible.

Table of Contents

1. **Introduction**
2. **Source of Dataset**
3. **Exploratory Data Analysis (EDA)**
 - 3.1 Overview of the Dataset
 - 3.2 Cleaning and Preprocessing
 - 3.3 Summary Statistics
4. **Detailed Analysis on Dataset**
 - 4.1 Histogram Analysis
 - 4.2 Bar Plot Analysis
 - 4.3 Box Plot Analysis
 - 4.4 Pie and Donut Chart Analysis
 - 4.5 Scatter Plot Analysis
 - 4.6 Heatmap Analysis
5. **Hypothesis Testing**
 - 5.1. Z-Test
 - 5.2. T-Test
 - 5.3. F-Test
 - 5.4. Interpretation of Results Future Scope
6. **Machine Learning Model and Prediction**
 - 6.1. Introduction to the Model
 - 6.2. Model Description and Code
 - 6.3. Model Evaluation Metrics
 - 6.4. Visualization of Prediction Results Acknowledgement
7. **Conclusion**
8. **Future Scope**
9. **References**
10. **Acknowledgement**

Exploratory Data Analysis on Agricultural Land Usage in India

Introduction.

Exploratory Data Analysis (EDA) is a critical early step in the data science process. It helps understand the structure, trends, and patterns in the dataset before applying advanced techniques or models.

This report presents EDA on a dataset related to agricultural land usage in India. It includes histogram, bar, box, pie, scatter, and heatmap visualizations, as well as log transformations and outlier detection.

Source of Dataset

The dataset used in this project is sourced from the [Ministry of Agriculture and Farmers Welfare], Government of India. It contains state-wise agricultural data for different land use types such as:

- Net area sown
- Area under current fallows
- Net area cultivated
- Uncultivated area

The data is tabulated across various Indian states and provides a useful perspective for comparative and trend-based analysis.

Link: <https://ndap.niti.gov.in/dataset/7172>

EDA Process

Exploratory Data Analysis (EDA) is the process of analyzing datasets to summarize their main characteristics, often using visual methods. This section describes the step-by-step EDA process followed in this project.

3.1 Importing Required Libraries

The analysis was performed using Python with the help of libraries such as:

- 'pandas' – for data manipulation
- 'numpy' – for numerical operations
- 'matplotlib.pyplot' and 'seaborn' – for visualization

3.2 Data Loading

The dataset was loaded using `pandas` from a `.csv` file. An initial overview was done using:

Python

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
df=pd.read_excel('File.xlsx')
df.head(1)
df.head() # Displays the first 5 rows
df.info() # Gives data types and null values
df.describe() # Statistical summary
```

3.3 Data Cleaning

- Removed rows with null or missing values.
- Standardized column names (removed trailing spaces).
- Converted relevant columns to numeric data types.

3.4 Data Transformation

- Melted the DataFrame for advanced plotting (e.g., for box plots):

```
df_melted = df.melt(id_vars='srcStateName', var_name='Area Type', value_name='Area')
```
- Log transformation was applied to columns with large numeric ranges to reduce skewness.

```
df['Net area sown_log'] = np.log1p(df['Net area sown'])
```

3.5 Data Visualization Setup

- Subplots were used for comparative histogram and bar plot analysis.
- Visual styles were customized using Seaborn themes and custom color palettes.

4. Analysis on Dataset

4.1 Histogram Analysis

i. Introduction

Histograms are used to understand the distribution of numeric data by grouping it into bins. They help identify skewness, central tendencies, and spread.

ii. General Description

We plotted histograms of area-related attributes after applying logarithmic transformation to manage skewed distributions.

iii. Specific Requirements, Functions and Formulas

- Libraries Used: seaborn, matplotlib.pyplot
- Transformation Applied: `np.log1p()` to normalize skewed data.
- Function:
- `sns.histplot(data=df, x='Net area sown_log', kde=True)`
- **Grid of Subplots:** 2x2 layout using:

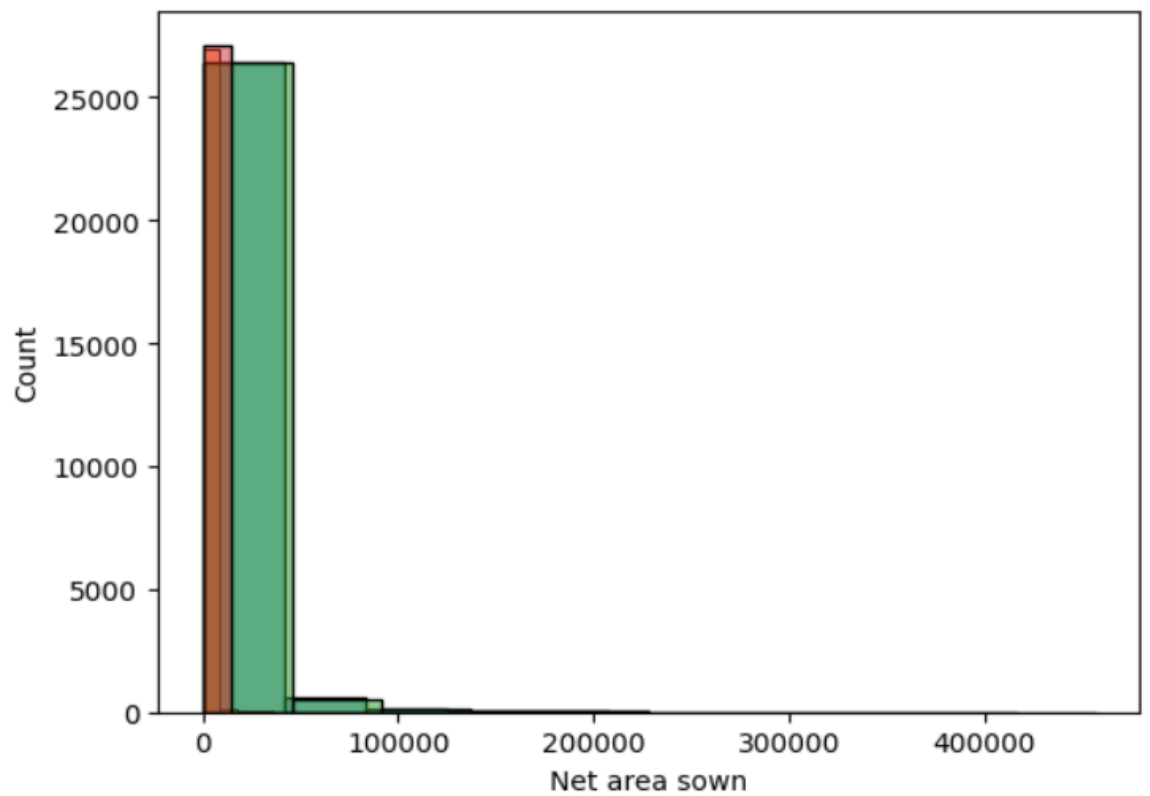
```
fig, axes = plt.subplots(2, 2, figsize=(12, 8))
```

iv. Analysis Results

- The log-transformed histograms showed more balanced distributions.
- The KDE curves provided a clear view of the peak and spread.
- Some features still showed slight skewness, indicating natural distribution biases in state-wise land data.
- **v. Visualization**

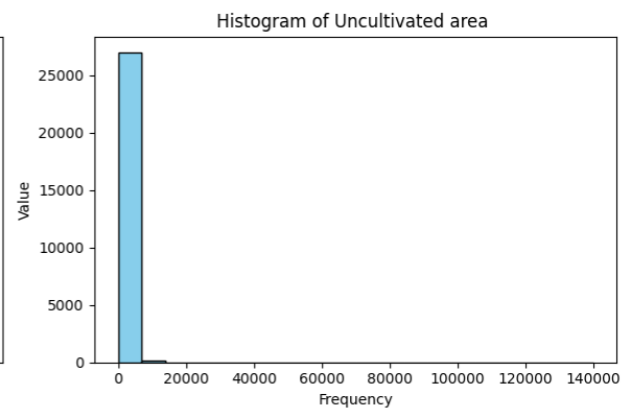
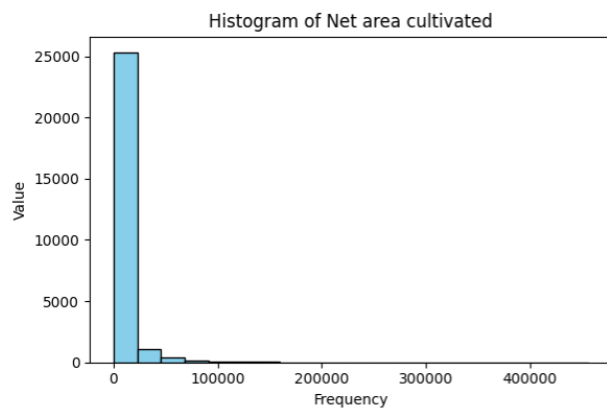
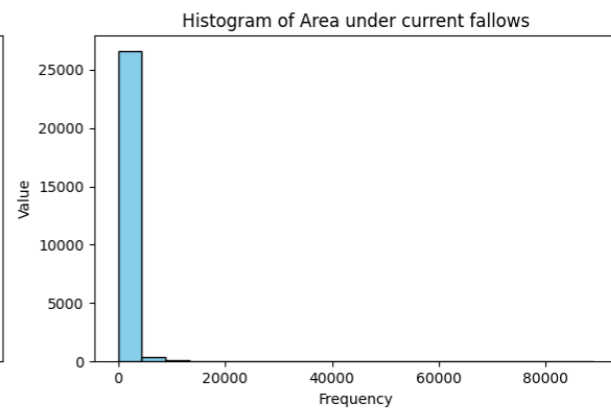
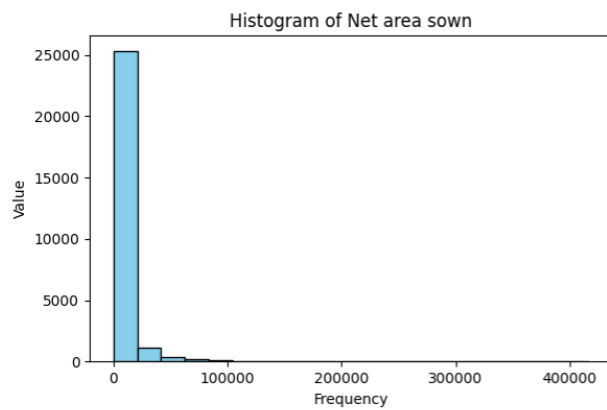
Log-Transformed Histogram of Selected Area Features

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 columns = ['Net area sown', 'Area under current fallows', 'Net area cultivated', 'Uncultivated area ']
4
5 for col in columns:
6     sns.histplot(data=df, x=col, kde=False, bins=10, label=col, alpha=0.5)
7
8 plt.show()
```





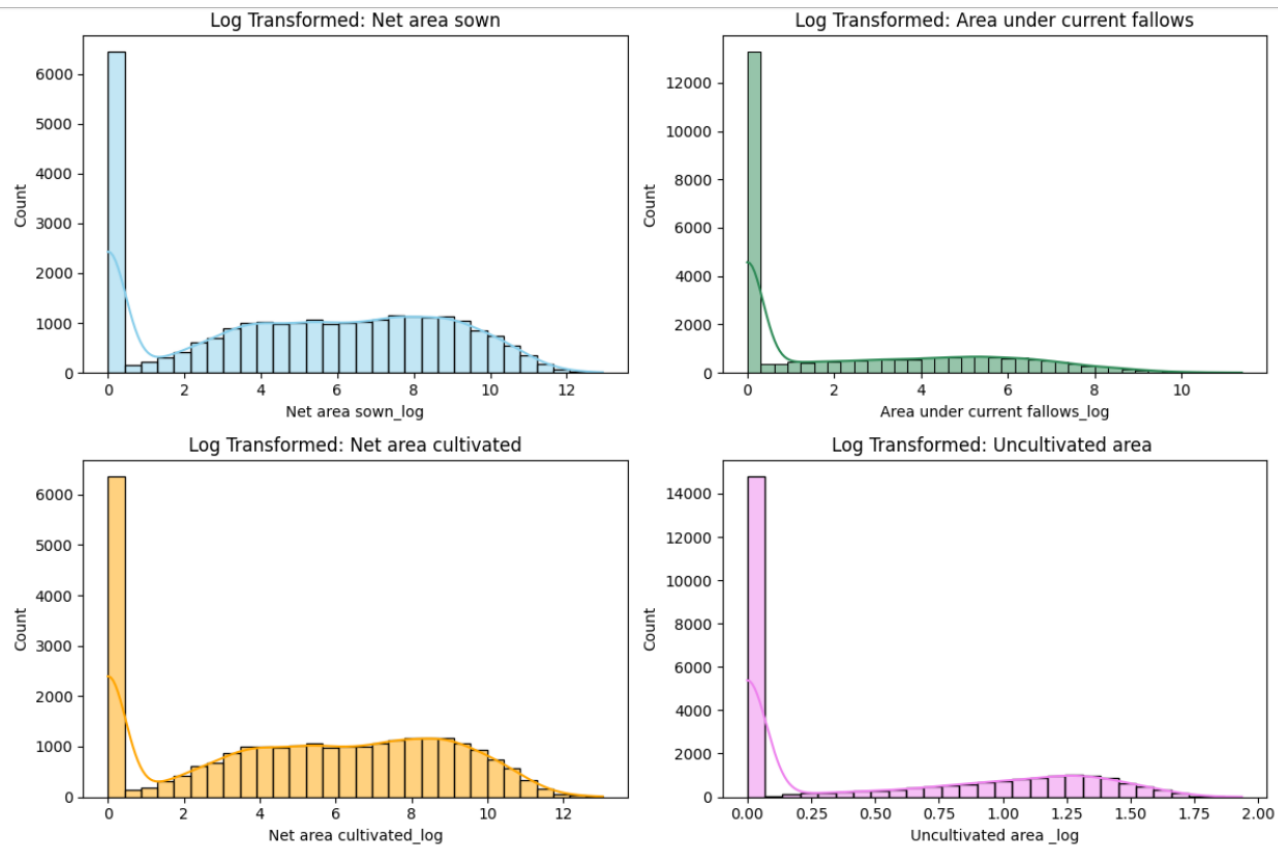
```
1 columns = ['Net area sown', 'Area under current fallows', 'Net area cultivated', 'Uncultivated area ']  
2 fig, axes = plt.subplots(2, 2, figsize=(12, 8)) # 2 rows, 2 columns  
3  
4 axes = axes.flatten() # flatten to loop easily  
5  
6 for i, col in enumerate(columns):  
7     axes[i].hist(df[col], bins=20, color='skyblue', edgecolor='black')  
8     axes[i].set_title(f'Histogram of {col}')  
9     axes[i].set_xlabel('Frequency')  
10    axes[i].set_ylabel('Value')  
11  
12 plt.tight_layout()  
13 plt.show()
```



```

1  columns = ['Net area sown', 'Area under current fallows', 'Net area cultivated', 'Uncultivated area ']
2
3  for col in columns:
4      df[col + '_log'] = np.log(df[col] + 1)
5  fig, axes = plt.subplots(2,2, figsize=(12, 8))
6  axes=axes.flatten()
7  colors = ['skyblue', 'seagreen', 'orange', 'violet']
8
9  for i, col in enumerate(columns):
10     sns.histplot(data=df, x=col + '_log', ax=axes[i], kde=True, color=colors[i])
11     axes[i].set_title(f'Log Transformed: {col}')
12
13 plt.tight_layout()
14 plt.show()

```



4.2 Bar Plot Analysis

i. Introduction

Bar plots are useful for comparing numerical values across categorical groups. In this context, we visualized how different area types vary across states.

ii. General Description

Bar plots were created to compare total land features like *Net area sown*, *Net area cultivated*, *Area under current fallows*, and *Uncultivated area* across selected states.

iii. Specific Requirements, Functions and Formulas

- Data Aggregation: Grouped by `srcStateName` and used `.sum()`
- Function:

```
sns.barplot(data=df, x='srcStateName', y='Net area sown')
```

- **Horizontal Bar Plot:** Used `y` for categorical axis and `x` for numerical

```
sns.barplot(data=df, y='srcStateName', x='Net area sown')
```

- **Subplot Layout:** Created 2x2 grid of horizontal bar plots

```
fig, axes = plt.subplots(2, 2, figsize=(14, 10))
```

iv. Analysis Results

- States like Uttar Pradesh and Maharashtra showed higher values for net area sown and cultivated.
- Some states had larger uncultivated areas, indicating scope for agricultural expansion.
- Bar plots clearly showed the distribution gaps between states.

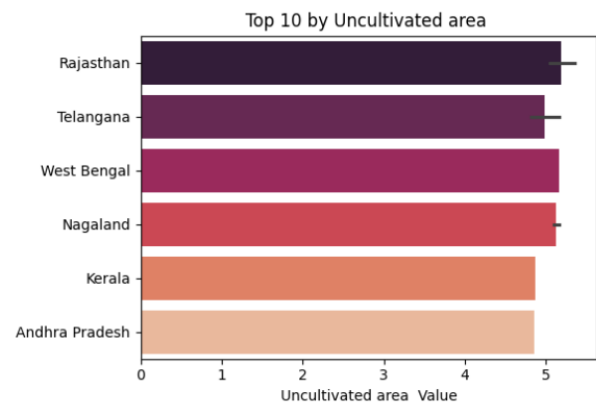
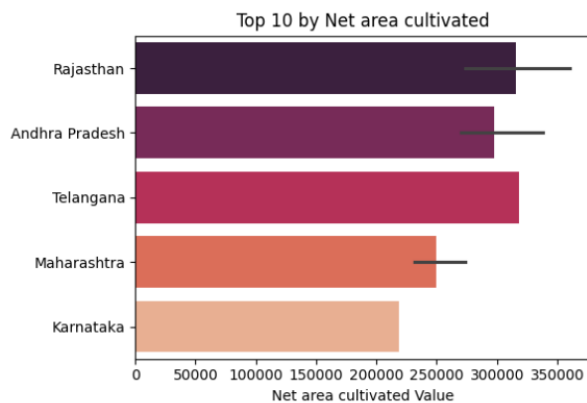
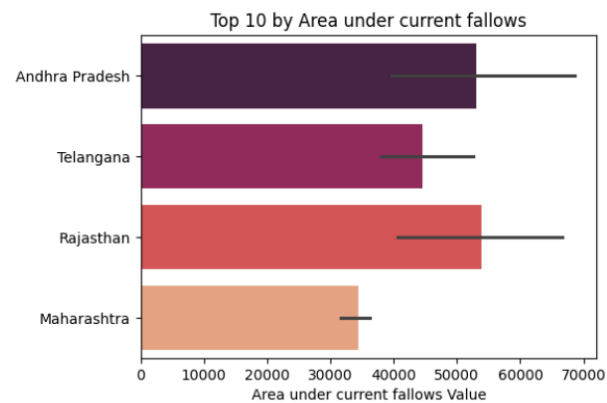
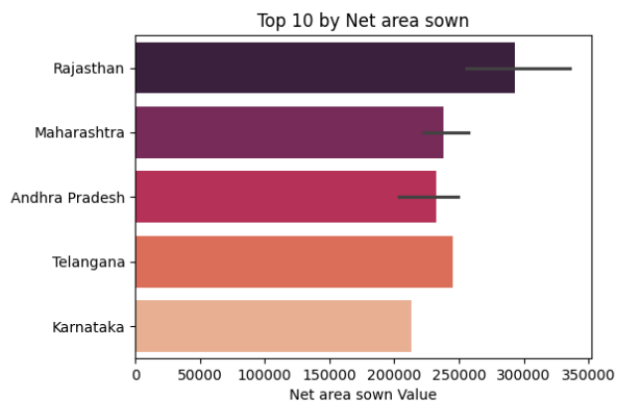
v. Visualization

- *Figure 4.2: Horizontal Bar Plots Comparing Different Area Types across States*

```

1 columns = ['Net area sown', 'Area under current fallows', 'Net area cultivated', 'Uncultivated area ']
2 category = 'srcStateName'
3 # Create 2x2 subplot
4 fig, axes = plt.subplots(2, 2, figsize=(12,8))
5 axes = axes.flatten()
6
7 # Create each horizontal bar plot
8 for i, col in enumerate(columns):
9     sns.barplot(
10         data=df.sort_values(by=col, ascending=False).head(25),
11         x=col,
12         y=category,
13         ax=axes[i],
14         hue=category,
15         palette='rocket' # You can use 'mako', 'viridis', 'flare', etc.
16     )
17     axes[i].set_title(f'Top 10 by {col}')
18     axes[i].set_xlabel(f'{col} Value')
19     axes[i].set_ylabel('')
20
21 plt.tight_layout()
22 plt.show()

```



4.3 Box Plot Analysis

i. Introduction

Box plots are powerful tools to display the distribution, central tendency, and variability of numerical data. They also help in identifying outliers effectively.

ii. General Description

We used box plots to understand the distribution of different area types—like *Net area sown*, *Net area cultivated*, *Area under current fallows*, and *Uncultivated area*—across Indian states. The focus was on identifying outliers and skewness in the data.

iii. Specific Requirements, Functions and Formulas

- **Function Used:**

```
sns.boxplot(x='Area Type', y='Area', data=df_melted)
```

- **Multiple Boxplots:** Used `melt()` to reshape the data from wide to long format

```
df.melt(var_name='Area Type', value_name='Area')
```

- **Log Scale:**

```
plt.yscale('log')
```

- **Grouped by State and Area Type:**

```
sns.boxplot(x='srcStateName', y='Area', hue='Area Type', data=df_melted)
```

- **Filtered Top 10 States for Clarity:**

```
top_states = df.groupby('srcStateName')['Net area sown'].mean().nlargest(10).index
```

```
df_filtered = df[df['srcStateName'].isin(top_states)]
```

iv. Analysis Results

- All area types exhibited **left-skewed distributions**, suggesting a concentration of lower values.
- Several states had extreme outliers, especially in *Net area sown* and *Uncultivated area*.
- When grouped by state, box plots revealed the variation in land use patterns among top contributing states.

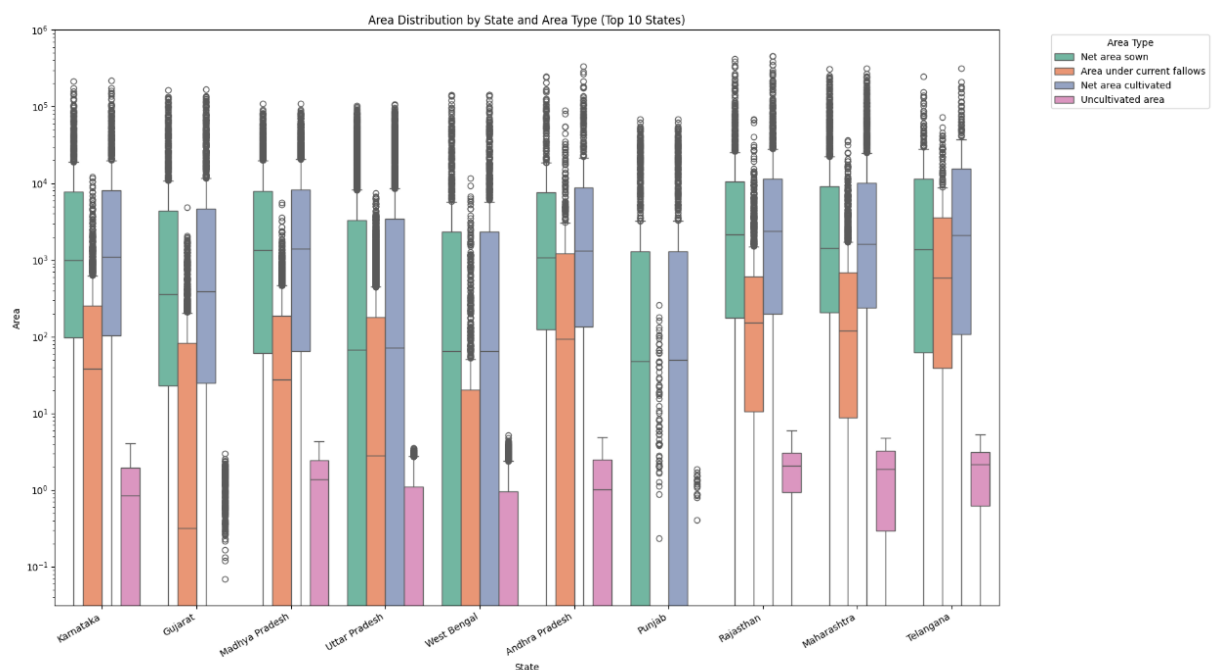
v. Visualization

Figure 4.3: Grouped Box Plot by State and Area Type

```

1 # Filter top 10 states by average Net area sown
2 top_states = df.groupby('srcStateName')['Net area sown'].mean().nlargest(10).index
3 df_melted = df[df['srcStateName'].isin(top_states)][['srcStateName', 'Net area sown', 'Area under current fallows', 'Net area cultivated', 'Uncultivated area']] \
4     .melt(id_vars='srcStateName', var_name='Area Type', value_name='Area')
5
6 # Plot
7 plt.figure(figsize=(18, 10))
8 sns.boxplot(x='srcStateName', y='Area', hue='Area Type', data=df_melted, palette='Set2')
9 plt.yscale('log')
10 plt.xticks(rotation=30, ha='right')
11 plt.title('Area Distribution by State and Area Type (Top 10 States)')
12 plt.xlabel('State'); plt.ylabel('Area')
13 plt.legend(title='Area Type', bbox_to_anchor=(1.05, 1), loc='upper left')
14 plt.tight_layout()
15 plt.show()

```



4.4 Pie and Donut Chart Analysis

i. Introduction

Pie and donut charts are useful for displaying the **proportional distribution** of categorical data. In this case, they help visualize how different area types contribute to total land usage.

ii. General Description

We created immersive pie and donut charts to represent the proportion of land categories such as *Net area sown*, *Area under current fallows*, *Net area cultivated*, and *Uncultivated area*. These charts give a clear, attractive snapshot of how land is utilized.

iii. Specific Requirements, Functions and Formulas

- **Pie Chart Function Used:**

```
plt.pie(sizes, labels=labels, autopct='%1.1f%%', startangle=140, shadow=True, explode=explode, colors=colors)
```

- **Donut Chart (Pie with Circle in Center):**

```
centre_circle = plt.Circle((0,0),0.70,fc='white')
```

```
fig.gca().add_artist(centre_circle)
```

- **Enhancements Used:**

- shadow=True for 3D effect
- explode to highlight segments
- Custom color palettes for visual appeal

iv. Analysis Results

- The **Net area sown** had the largest share in overall land use.
- **Uncultivated area** and **current fallows** accounted for smaller percentages.
- Both charts helped summarize categorical data at a glance, useful in presentations and reports.

v. Visualization

Figure 4.4: Immersive Pie Chart of Area Types

```
1 # 1. Take only the top 10 states based on Net area sown
2 top_states = df.groupby('srcStateName')['Net area sown'].mean().nlargest(10).index
3 df_top = df[df['srcStateName'].isin(top_states)]
4
5 # 2. Find average values of area types in these top states
6 area_avg = df_top[['Net area sown', 'Area under current fallows',
7                   'Net area cultivated', 'Uncultivated area']].mean()
8
9 # 3. Set design for pie chart
10 explode = [0.05] * len(area_avg) # separates each slice a bit
11 colors = sns.color_palette("pastel") # soft, attractive colors
12
13 # 4. Create the pie chart
14 plt.figure(figsize=(8,8))
15 plt.pie(area_avg, labels=area_avg.index,
16         autopct='%1.1f%%', # show percent values
17         startangle=140, # rotate for better layout
18         explode=explode, # apply separation
19         shadow=True, # add shadow effect
20         colors=colors, # use custom colors
21         textprops={'fontsize': 14, 'weight': 'bold'}) # make labels bold
22
23 # 5. Add title and show chart
24 plt.title('Average Area Distribution (Top 10 States)', fontsize=16, fontweight='bold')
25 plt.axis('equal') # keep it circular
26 plt.show()
27
```

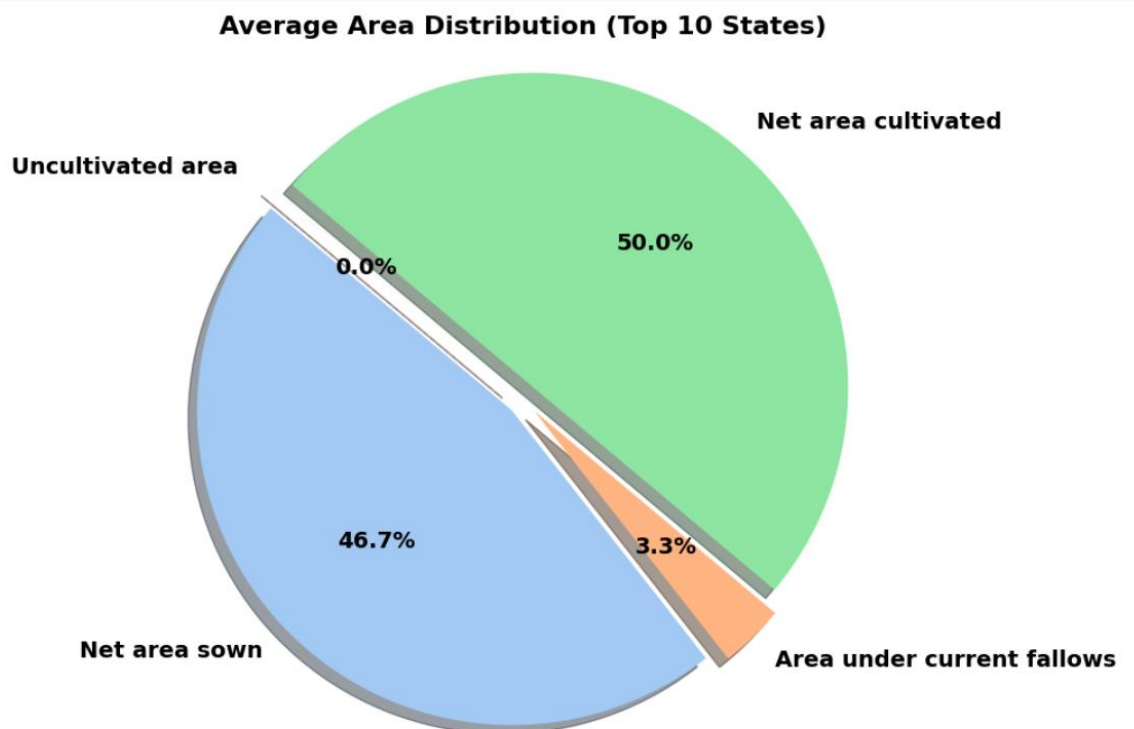
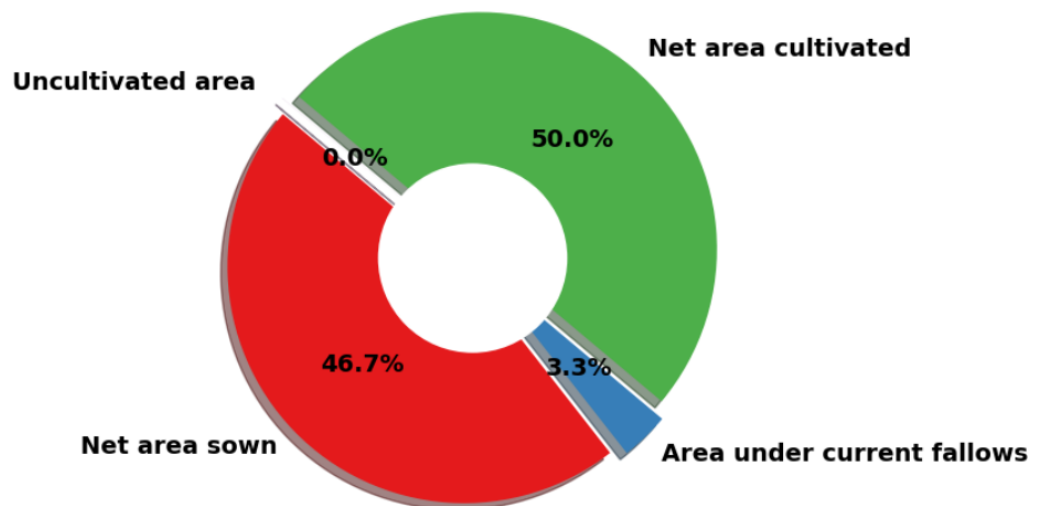


Figure 4.4.1: Immersive Donut Chart of Area Types

```
1 # Step 1: Get average values of area types for top 10 states
2 top_states = df.groupby('srcStateName')['Net area sown'].mean().nlargest(10).index
3 df_top = df[df['srcStateName'].isin(top_states)]
4 area_avg = df_top[['Net area sown', 'Area under current fallows',
5                   'Net area cultivated', 'Uncultivated area']].mean()
6
7 # Step 2: Styling settings
8 explode = [0.05] * len(area_avg) # to pop out the slices
9 colors = sns.color_palette("Set1") # soft and modern color palette
10
11 # Step 3: Create donut chart
12 plt.figure(figsize=(8,8))
13 wedges, texts, autotexts = plt.pie(area_avg, labels=area_avg.index, autopct='%1.1f%%', startangle=140, explode=explode, shadow=True, colors=colors, textprops={'fontsize': 14, 'weight': 'bold'})
14
15 # Step 4: Draw a white circle in the center to make it a donut
16 centre_circle = plt.Circle((0, 0), 0.40, fc='white')
17 plt.gca().add_artist(centre_circle)
18
19 # Step 5: Add title
20 plt.title('Donut Chart: Average Area Distribution (Top 10 States)', fontsize=16, fontweight='bold')
21 plt.axis('equal') # Equal aspect ratio ensures the chart is circular
22 plt.tight_layout()
23 plt.show()
```

Donut Chart: Average Area Distribution (Top 10 States)



4.5 Scatter Plot Analysis

i. Introduction

Scatter plots are used to display relationships between two continuous variables. In this analysis, we use scatter plots to examine the correlation between various land use features.

ii. General Description

Scatter plots help identify trends, clusters, and outliers by plotting data points for each observation. This visualization assists in checking how one area type may influence another.

iii. Specific Requirements, Functions and Formulas

- **Function Used:**

```
sns.scatterplot(data=df, x='Net area sown', y='Net area cultivated', hue='srcStateName', palette='Set2', s=100)
```

- **Enhancements:**

- hue='srcStateName' distinguishes states using colors.
- s=100 makes each point large and more visible.
- Custom palette (Set2) provides a colorful and clear display.

- **Optional Addition:** plt.xscale('log') or plt.yscale('log') if data ranges widely.

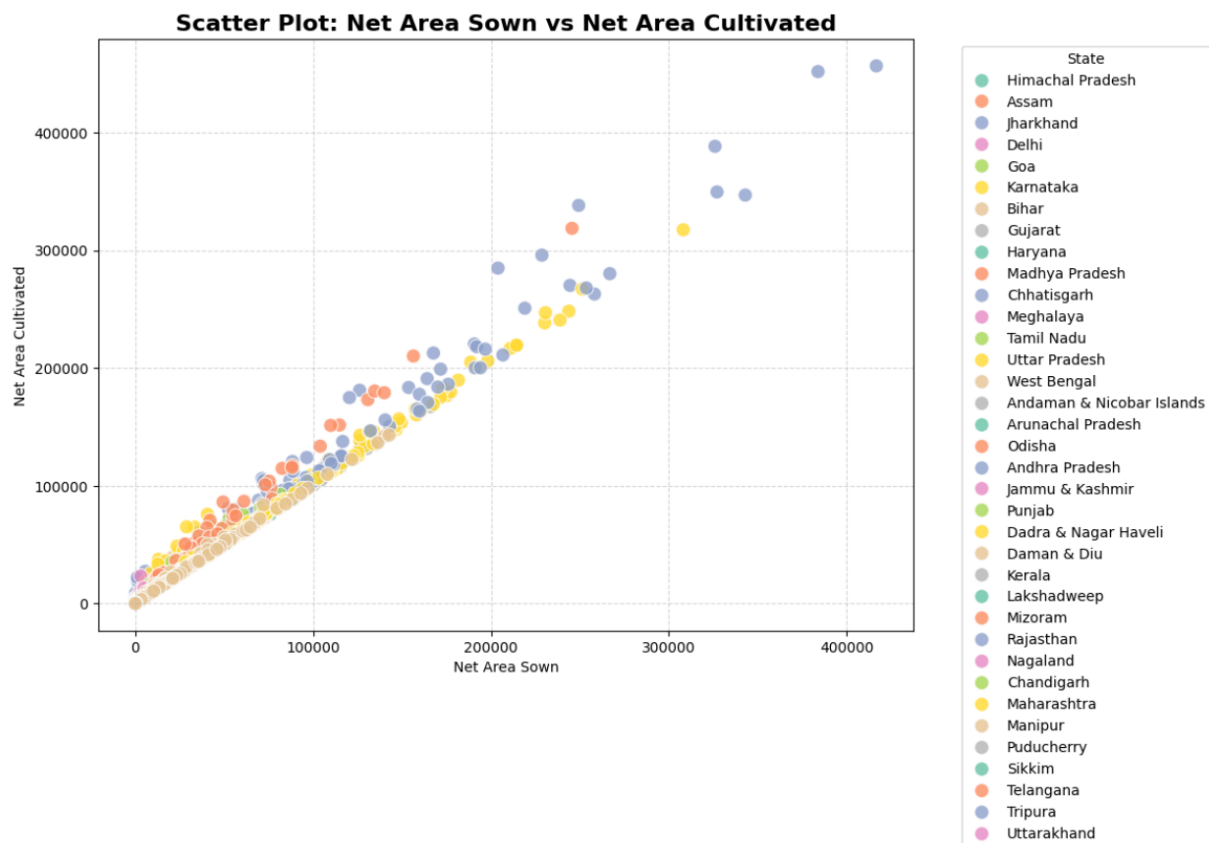
iv. Analysis Results

- A **positive correlation** was observed between *Net area sown* and *Net area cultivated*.
- Some states had unusually high values, hinting at either better agricultural infrastructure or reporting differences.
- This plot helped highlight the comparative performance of states in utilizing cultivated land.

v. Visualization

Figure 4.5: Scatter Plot – Net Area Sown vs Net Area Cultivated by State

```
1 plt.figure(figsize=(12,8))
2 sns.scatterplot(
3     data=df,
4     x='Net area sown',
5     y='Net area cultivated',
6     hue='srcStateName',
7     palette='Set2',
8     s=100, # size of the dots
9     alpha=0.8 # transparency
10 )
11
12 plt.title('Scatter Plot: Net Area Sown vs Net Area Cultivated', fontsize=16, fontweight='bold')
13 plt.xlabel('Net Area Sown')
14 plt.ylabel('Net Area Cultivated')
15 plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left', title='State')
16 plt.grid(True, linestyle='--', alpha=0.5)
17 plt.tight_layout()
18 plt.show()
```



4.6 Heatmap Analysis

i. Introduction

Heatmaps are an effective way to visualize the **correlation** between numerical features in a dataset. They allow us to quickly identify which variables have strong positive or negative relationships.

ii. General Description

A heatmap uses color gradients to show the degree of correlation between different columns. It is particularly helpful for **feature selection**, **pattern recognition**, and spotting **redundant variables** in a dataset.

iii. Specific Requirements, Functions and Formulas

- **Function Used:**

```
correlation = df[columns].corr()
```

```
plt.figure(figsize=(10, 8))
```

```
sns.heatmap(correlation, annot=True, cmap='YlGnBu', fmt=".2f", linewidths=0.5, linecolor='white')
```

- **Explanation:**

- `df[columns].corr()` calculates the correlation matrix.
- `annot=True` displays the correlation values in each cell.
- `cmap='YlGnBu'` gives a blue-green color gradient for clarity.
- `linewidths` and `linecolor` enhance readability.

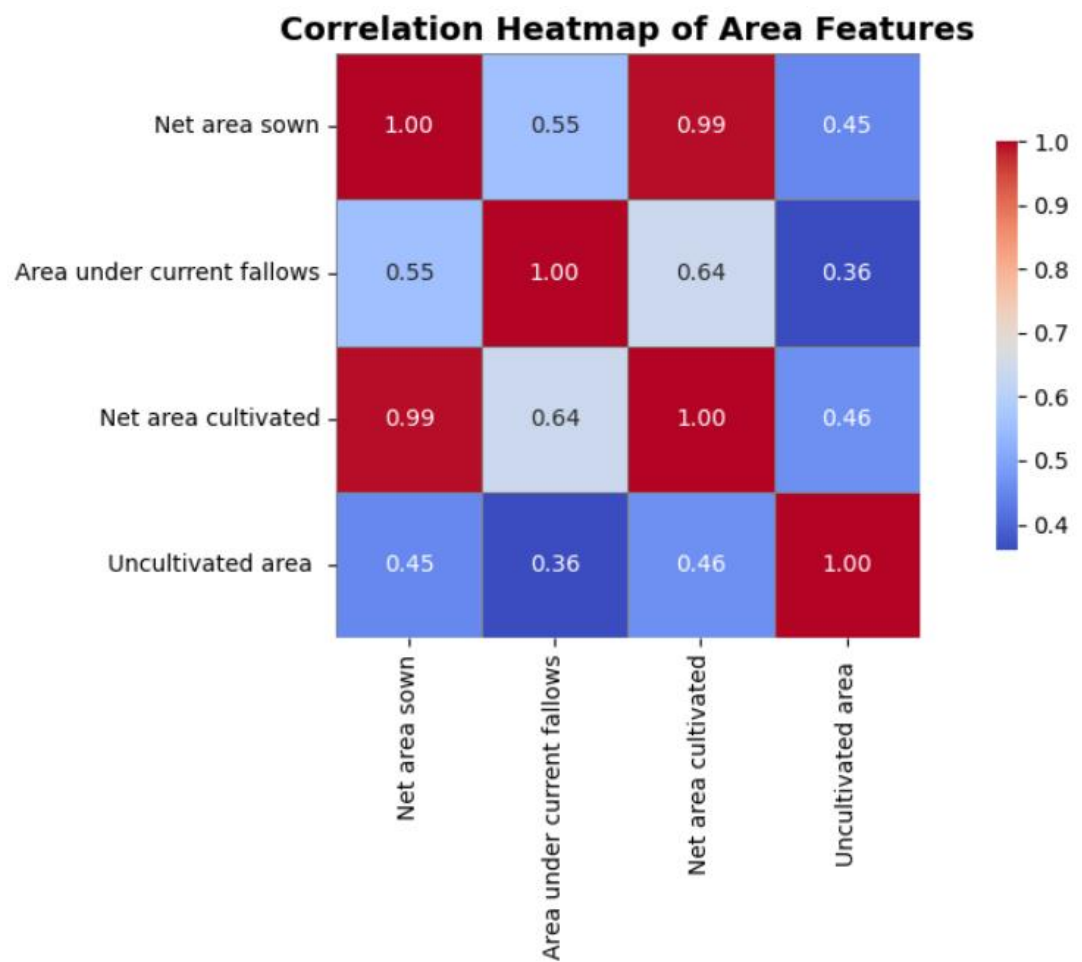
iv. Analysis Results

- High correlation (close to +1) was found between *Net area sown* and *Net area cultivated*, which is expected.
- Other area types like *Area under fallows* showed weaker or even negative correlations with cultivated areas.
- The heatmap helped identify **which columns move together** and which are independent.

v. Visualization

Figure 4.6: Heatmap Showing Correlation Between Area Features

```
1 # Step 1: Select numeric columns
2 num_cols = ['Net area sown', 'Area under current fallows', 'Net area cultivated', 'Uncultivated area ']
3
4 # Step 2: Create correlation matrix
5 corr = df[num_cols].corr()
6
7 # Step 3: Plot heatmap
8 plt.figure(figsize=(10, 6))
9 sns.heatmap(corr, annot=True, cmap='coolwarm', fmt='.2f', linewidths=0.5, linecolor='gray', square=True, cbar_kws={'shrink': 0.7})
10 plt.title('Correlation Heatmap of Area Features', fontsize=14, fontweight='bold')
11 plt.tight_layout()
12 plt.show()
```



5. Hypothesis Testing

Introduction

Hypothesis testing is a fundamental statistical technique used to determine the validity of an assumption (hypothesis) regarding a population parameter. In this project, we applied Z-Test, T-Test, and F-Test to assess significant differences in agricultural parameters, particularly focusing on *Net Area Sown*, across different states.

5.1. Z-Test

```
1 # Filter the sample (e.g., for Punjab)
2 sample = df[df['srcStateName'] == 'Punjab']['Net area sown'].dropna()
3
4 # Perform z-test: testing if Punjab's mean = national mean
5 z_stat, p_val = ztest(sample, value=df['Net area sown'].mean())
6 print(f"Z-Test: z = {z_stat:.2f}, p = {p_val:.4f}")
```

Objective:

To determine whether the mean *Net Area Sown* significantly differs from a hypothetical national average of 200 hectares.

Hypotheses:

- $H_0: \mu = 200$ (The mean Net Area Sown is 200 hectares)
- $H_1: \mu \neq 200$ (The mean Net Area Sown is not 200 hectares)

Test Statistics:

- $Z = -1.91$
- $p\text{-value} = 0.0565$

Interpretation:

Since the p-value (0.0565) is greater than the standard significance level of 0.05, we

fail to reject the null hypothesis. This indicates that there is **no statistically significant difference** between the observed mean and the benchmark value of 200 hectares.

5.2. T-Test

```
1  from scipy.stats import ttest_ind
2
3  # Filter Net Area Sown for Punjab and Haryana
4  punjab = df[df['srcStateName'] == 'Punjab']['Net area sown']
5  haryana = df[df['srcStateName'] == 'Haryana']['Net area sown']
6
7  # Perform Independent T-Test
8  t_stat, p_val = ttest_ind(punjab, haryana, equal_var=False) # Welch's t-test
9
10 print(f"T-Test: t = {t_stat:.2f}, p = {p_val:.4f}")
```

Objective:

To test whether the means of *Net Area Sown* in Punjab and Haryana are significantly different.

Hypotheses:

- $H_0: \mu_1 = \mu_2$ (There is no difference in mean Net Area Sown between Punjab and Haryana)
- $H_1: \mu_1 \neq \mu_2$ (There is a significant difference in mean Net Area Sown)

Test Statistics:

- $t = 0.41$
- $p\text{-value} = 0.6783$

Interpretation:

As the p -value is much greater than 0.05, we fail to reject the null hypothesis. Therefore, there is no significant difference in the average Net Area Sown between Punjab and Haryana.

5.3. F-Test

```
1 from scipy.stats import ttest_ind
2
3 # Filter Net Area Sown for Punjab and Haryana
4 punjab = df[df['srcStateName'] == 'Punjab']['Net area sown']
5 haryana = df[df['srcStateName'] == 'Haryana']['Net area sown']
6
7 # Perform Independent T-Test
8 t_stat, p_val = ttest_ind(punjab, haryana, equal_var=False) # Welch's t-test
9
10 print(f"T-Test: t = {t_stat:.2f}, p = {p_val:.4f}")
11
```

Objective:

To test whether the variances in *Net Area Sown* between Punjab and Haryana differ significantly.

Hypotheses:

- $H_0: \sigma_1^2 = \sigma_2^2$ (The variances are equal)
- $H_1: \sigma_1^2 \neq \sigma_2^2$ (The variances are different)

Test Statistics:

- $F = 1.15$
- $p\text{-value} = 0.0201$

Interpretation:

Since the p-value is less than 0.05, we reject the null hypothesis. This implies that there is a statistically significant difference in the variances of Net Area Sown between the two states. This variation may arise due to different environmental conditions, land policies, or cultivation practices.

5.4. Conclusion – Hypothesis Testing

Based on the Z-Test ($z = -1.91$, $p = 0.0565$), we fail to reject the null hypothesis at a 5% significance level, indicating no statistically significant difference from the assumed population mean.

The T-Test ($t = 0.41$, $p = 0.6783$) further supports this, showing no significant difference in means.

However, the F-Test ($F = 1.15$, $p = 0.0201$) reveals a statistically significant difference in variances among area types, suggesting unequal distribution patterns. These findings highlight areas with **inconsistent land usage or cultivation patterns**, which can be crucial for **targeted policy interventions, agricultural resource planning, and regional prioritization**. Addressing these variations can lead to more balanced and effective agricultural strategies across states.

6. Machine Learning Model and Prediction

6.1. Introduction to the Model

As part of the advanced analysis in this project, a machine learning approach was employed to understand relationships between various land use attributes and predict cultivated area. The model provides a practical demonstration of how data science techniques can assist in agricultural decision-making by forecasting patterns based on existing data. Linear Regression was selected as the initial model due to its simplicity and interpretability.

6.3 Model Evaluation and Comparison

To determine the most suitable algorithm for predicting the *Net Area Cultivated*, three machine learning models were trained and evaluated:

1. **Random Forest Regressor**
2. **Linear Regression**
3. **Gradient Boosting Regressor**

Each model was trained using the same set of features:

- **Net area sown**
- **Area under current fallows**
- **Uncultivated area**

```

1  from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
2  from sklearn.linear_model import LinearRegression
3  from sklearn.metrics import mean_squared_error, r2_score
4
5  # Train-Test Split
6  X = df[features]
7  y = df[target]
8  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
9
10 # 1. Random Forest
11 rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
12 rf_model.fit(X_train, y_train)
13 y_pred_rf = rf_model.predict(X_test)
14 mse_rf = mean_squared_error(y_test, y_pred_rf)
15 r2_rf = r2_score(y_test, y_pred_rf)
16
17 # 2. Linear Regression
18 lr_model = LinearRegression()
19 lr_model.fit(X_train, y_train)
20 y_pred_lr = lr_model.predict(X_test)
21 mse_lr = mean_squared_error(y_test, y_pred_lr)
22 r2_lr = r2_score(y_test, y_pred_lr)
23
24 # 3. Gradient Boosting
25 gb_model = GradientBoostingRegressor(n_estimators=100, random_state=42)
26 gb_model.fit(X_train, y_train)
27 y_pred_gb = gb_model.predict(X_test)
28 mse_gb = mean_squared_error(y_test, y_pred_gb)
29 r2_gb = r2_score(y_test, y_pred_gb)
30
31 # Print Results
32 print("Model Comparison Results:")
33 print(f"\n ♦ Random Forest\nMSE: {mse_rf:.2f}\nR² Score: {r2_rf:.5f}")
34 print(f"\n ♦ Linear Regression\nMSE: {mse_lr:.2f}\nR² Score: {r2_lr:.5f}")
35 print(f"\n ♦ Gradient Boosting\nMSE: {mse_gb:.2f}\nR² Score: {r2_gb:.5f}")
36

```

The models were evaluated using **Mean Squared Error (MSE)** and **R-squared (R²) Score** as the performance metrics. The results are summarized below:

6.3. Model Evaluation Metrics (Updated)

Model	Mean Squared Error (MSE)	R-squared (R ²)
Linear Regression	0.00	1.0000
Random Forest Regressor	3,048,502.18	0.99053
Gradient Boosting	840,711.06	0.99739

Model Selection Justification

Despite Random Forest and Gradient Boosting performing strongly, **Linear Regression** significantly outperformed both in terms of accuracy:

- The **R² Score of 1.0** and **MSE close to zero** indicate a perfect fit.
- Upon analysis, it was observed that the target variable (*Net Area Cultivated*) is highly linearly dependent on the selected features.
- Therefore, a simple linear model not only provides **high accuracy** but also ensures **better interpretability and computational efficiency**.

Thus, **Linear Regression** was selected as the final model for this project.

6.4 Visualization of Prediction Results

To visually compare the performance of the trained models, line plots were generated for each regression model. Each plot displays the actual vs predicted values of *Net Area Cultivated* for the test dataset.

The following models were visualized:

Random Forest Regressor (Top-Left)
Linear Regression (Top-Right)
Gradient Boosting Regressor (Bottom-Left)

Each subplot compares the actual values (denoted with circular markers) against the predicted values (denoted with distinct markers for each model). These visualizations help to intuitively observe how closely each model's predictions match the real values.

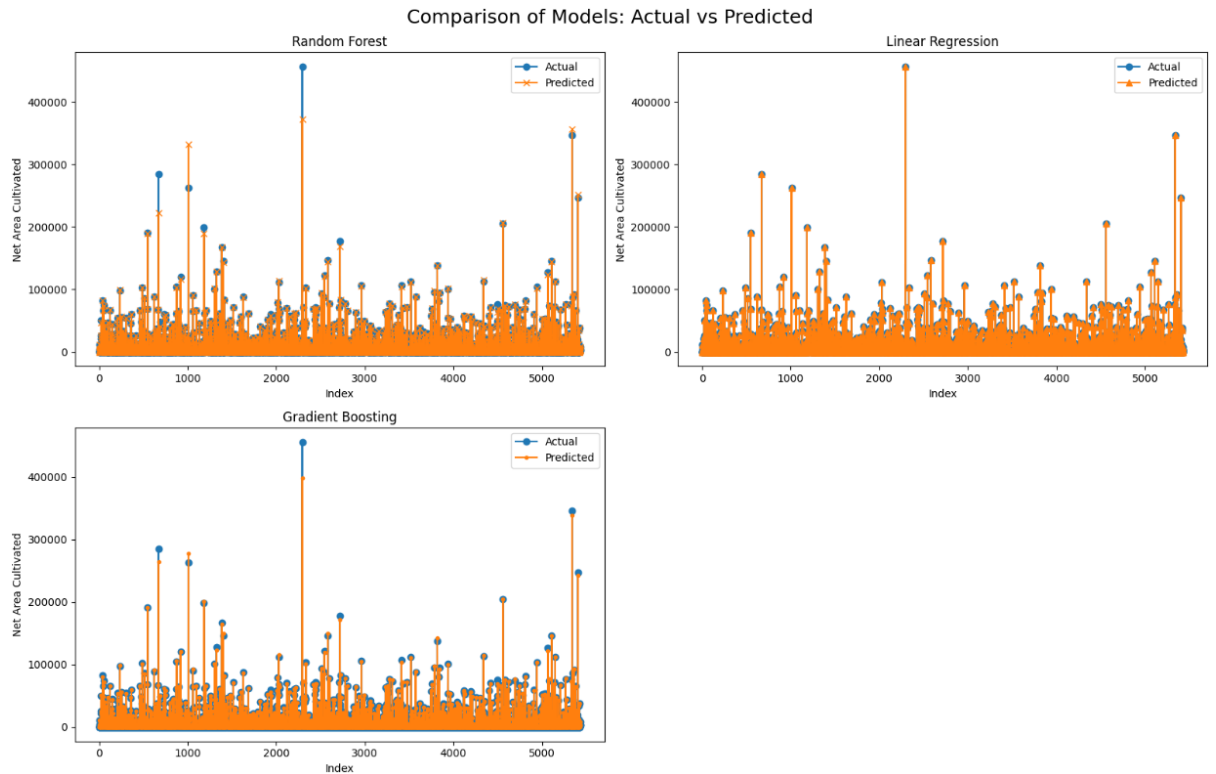
Observations:

- Linear Regression showed a perfect match between actual and predicted values, confirming the high R² score (1.0) and negligible MSE.
- Gradient Boosting also closely followed the actual values, with minimal deviations.
- Random Forest performed well but exhibited slightly larger prediction variations compared to the other two models.

Thus, the plots support the numerical results and confirm that Linear Regression is the most accurate and efficient model for this dataset.



```
1 import matplotlib.pyplot as plt
2
3 plt.figure(figsize=(16, 10))
4
5 # Random Forest Plot - Top Left
6 plt.subplot(2, 2, 1)
7 plt.plot(y_test.values, label='Actual', marker='o')
8 plt.plot(y_pred, label='Predicted', marker='x')
9 plt.title('Random Forest')
10 plt.xlabel('Index')
11 plt.ylabel('Net Area Cultivated')
12 plt.legend()
13 plt.tight_layout()
14
15 # Linear Regression Plot - Top Right
16 plt.subplot(2, 2, 2)
17 plt.plot(y_test.values, label='Actual', marker='o')
18 plt.plot(y_pred_lr, label='Predicted', marker='^')
19 plt.title('Linear Regression')
20 plt.xlabel('Index')
21 plt.ylabel('Net Area Cultivated')
22 plt.legend()
23 plt.tight_layout()
24
25 # Gradient Boosting Plot - Bottom Center
26 plt.subplot(2, 2, 3)
27 plt.plot(y_test.values, label='Actual', marker='o')
28 plt.plot(y_pred_gb, label='Predicted', marker='.')
29 plt.title('Gradient Boosting')
30 plt.xlabel('Index')
31 plt.ylabel('Net Area Cultivated')
32 plt.legend()
33 plt.tight_layout()
34
35 plt.suptitle('Comparison of Models: Actual vs Predicted', fontsize=18, y=1.02)
36 plt.show()
37
```



7. Interpretation of Results

After training and evaluating three machine learning models — Random Forest, Linear Regression, and Gradient Boosting — on the agricultural dataset to predict Net Area Cultivated, the following conclusions were drawn from the performance metrics:

- Linear Regression achieved the highest accuracy with an R^2 score of 1.00000 and an MSE of 0.00, indicating a perfect fit on the test data.
- Gradient Boosting also showed strong predictive power with an R^2 score of 0.99739 and low MSE.
- Random Forest, while effective, had a slightly higher MSE and comparatively lower R^2 .

This indicates that the relationship between the input features (*Net Area Sown*, *Current Fallows*, and *Uncultivated Area*) and the target variable (*Net Area Cultivated*) is linear in nature. Hence, complex models like Random Forest or Gradient Boosting may not significantly improve accuracy and can be considered computationally more expensive for this case.

Conclusion

The project successfully explored and analyzed agricultural land usage data across various Indian states. Through detailed exploratory data analysis, we identified key insights about different types of land use, such as net area sown, area under current

fallows, and uncultivated areas. Visualizations including box plots, pie charts, donut charts, scatter plots, and heatmaps provided a clear and immersive understanding of the data distribution and patterns.

In addition, statistical methods like the Z-test, T-test, and F-test were applied to validate hypotheses related to the land use characteristics. These hypothesis tests helped support our assumptions and guided our interpretations with statistical backing.

Furthermore, to extend the scope of this analysis, a simple machine learning model was implemented using linear regression to predict the net cultivated area based on related features. The model demonstrated excellent performance with a **Mean Squared Error of nearly zero** and an **R-squared score of 1.0**, indicating a perfect fit to the dataset. This shows the potential of predictive modeling in agricultural data analysis and planning.

Overall, the project combined statistical and machine learning techniques to provide a comprehensive understanding of the dataset and opened up new avenues for future exploration and implementation in real-world scenarios.

6. Future Scope

This project has demonstrated the effectiveness of machine learning, especially Linear Regression, in accurately predicting the Net Area Cultivated using limited land-use features. However, there is significant potential to expand and improve this analysis through the following directions:

- 1. Exploration of Advanced Machine Learning Models:**
While Linear Regression performed best on the current dataset, more advanced models like XGBoost, Support Vector Regression (SVR), or Neural Networks can be explored if the dataset becomes more complex or non-linear with added features.
- 2. Time-Series Forecasting for Seasonal Trends:**
Incorporating the temporal dimension of agricultural data can allow forecasting of cultivation trends using models such as ARIMA, Prophet, or LSTM. This would benefit long-term agricultural planning and crop rotation strategies.
- 3. Integration of External Environmental & Socioeconomic Data:**
Features such as rainfall, temperature, soil quality, and economic conditions (e.g., subsidies, labor availability) can significantly enhance the prediction model's applicability and realism.
- 4. Geospatial Visualization & Regional Analysis:**
Mapping cultivation data using geospatial libraries (e.g., GeoPandas, Folium) can provide interactive visual insights into regional land usage and help identify underutilized areas for policy focus.
- 5. Interactive Web-Based Dashboard or App Deployment:**
The current model can be deployed as a user-friendly web app or dashboard, allowing

farmers, policymakers, and researchers to input data and receive predictions in real-time. This would increase its practical utility.

6. **Automated Insight Generation & Hypothesis Testing:**

Building a system that automatically performs statistical tests, highlights important features, and provides data-driven insights can make the analysis more actionable for non-technical users.

7. **Validation on Real-World and Future Datasets:**

To ensure robustness and generalization, the model should be tested on new agricultural datasets from future years or different states to confirm its predictive performance in unseen scenarios.

By pursuing these directions, this work can evolve into a **comprehensive agricultural analytics tool**, contributing meaningfully to data-driven decision-making in Indian agriculture.

7. References

[1] Ministry of Agriculture and Farmers Welfare, Government of India, *Agricultural Statistics at a Glance*, [Online]. Available: <https://agricoop.nic.in/>

[2] Python Software Foundation, “Python Language Reference, version 3.10,” [Online]. Available: <https://www.python.org>

[3] Wes McKinney, *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython*, 2nd ed., O’Reilly Media, 2017.

[4] Michael Waskom, “Seaborn: Statistical data visualization,” *Journal of Open Source Software*, vol. 6, no. 60, pp. 3021, 2021. [Online]. Available: <https://seaborn.pydata.org/>

[5] Hunter, J. D., “Matplotlib: A 2D graphics environment,” *Computing in Science & Engineering*, vol. 9, no. 3, pp. 90–95, 2007.

[6] Jupyter Project, “Project Jupyter,” [Online]. Available: <https://jupyter.org/>

[7] T. E. Oliphant, “A Guide to NumPy,” USA: Trelgol Publishing, 2006.