

# Comprehensive Guide to pandas with the Iris Dataset

Sayan Chaki

March 12, 2025

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Loading the Dataset</b>	<b>3</b>
2.1	Using scikit-learn . . . . .	3
2.2	Using seaborn . . . . .	3
<b>3</b>	<b>Basic Exploratory Data Analysis</b>	<b>4</b>
3.1	Dataset Information . . . . .	4
3.2	Summary Statistics . . . . .	4
3.3	Data Distribution by Species . . . . .	5
<b>4</b>	<b>Data Manipulation Techniques</b>	<b>5</b>
4.1	Selecting and Filtering Data . . . . .	5
4.2	Adding New Columns . . . . .	6
4.3	Grouping and Aggregation . . . . .	6
<b>5</b>	<b>Statistical Analysis</b>	<b>6</b>
5.1	Correlation Analysis . . . . .	6
5.2	Descriptive Statistics by Group . . . . .	7
<b>6</b>	<b>Data Visualization with pandas</b>	<b>7</b>
6.1	Histograms . . . . .	7
6.2	Box Plots . . . . .	8
6.3	Scatter Plots . . . . .	9
<b>7</b>	<b>Advanced pandas Techniques</b>	<b>10</b>
7.1	Pivot Tables . . . . .	10
7.2	MultiIndex and Advanced Indexing . . . . .	11
7.3	Applying Functions with apply() . . . . .	11
<b>8</b>	<b>Exporting Data for LaTeX</b>	<b>12</b>
8.1	Converting DataFrames to LaTeX Tables . . . . .	12
8.2	Customizing LaTeX Output . . . . .	12

<b>9 Working with Missing Data</b>	<b>12</b>
<b>10 Merging and Joining DataFrames</b>	<b>13</b>
<b>11 Advanced Filtering and Selection</b>	<b>13</b>
<b>12 Pandas with Iris Dataset: Case Studies</b>	<b>14</b>
12.1 Case Study 1: Species Classification . . . . .	14
12.2 Case Study 2: Feature Engineering . . . . .	14
<b>13 Conclusion</b>	<b>15</b>
<b>A Complete Code Example</b>	<b>15</b>

# 1 Introduction

This guide provides a comprehensive overview of using pandas with the Iris dataset. The Iris dataset is a classic dataset in statistics and machine learning, introduced by the British statistician and biologist Ronald Fisher in 1936. It consists of 150 samples from three species of Iris (Iris setosa, Iris virginica, and Iris versicolor), with four features measured for each sample: the length and width of the sepals and petals.

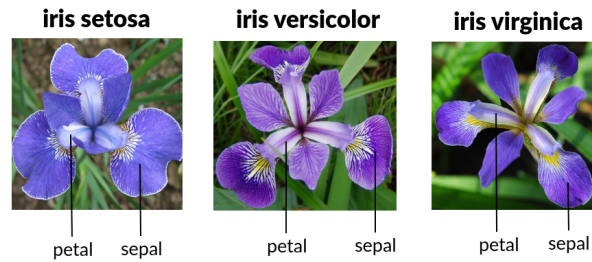


Figure 1: The three species of Iris in the dataset

## 2 Loading the Dataset

The first step in working with the Iris dataset is to load it into a pandas DataFrame. This can be done in several ways:

### 2.1 Using scikit-learn

```
1 from sklearn import datasets
2 import pandas as pd
3
4 # Load the Iris dataset
5 iris = datasets.load_iris()
6
7 # Create a DataFrame
8 df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
9 df['species'] = pd.Categorical.from_codes(iris.target, iris.target_names)
10
11 # Display the first few rows
12 print(df.head())
```

Listing 1: Loading Iris dataset from scikit-learn

### 2.2 Using seaborn

```
1 import seaborn as sns
2 import pandas as pd
3
4 # Load the Iris dataset
```

```

5 iris = sns.load_dataset('iris')
6
7 # Display the first few rows
8 print(iris.head())

```

Listing 2: Loading Iris dataset from seaborn

## 3 Basic Exploratory Data Analysis

### 3.1 Dataset Information

Once the dataset is loaded, we can obtain basic information about it:

```

1 # Basic information about the dataset
2 print(df.info())
3
4 # Summary statistics
5 print(df.describe())
6
7 # Check for missing values
8 print(df.isnull().sum())

```

Listing 3: Getting dataset information

The output of `df.info()` would look like:

```

1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 150 entries, 0 to 149
3 Data columns (total 5 columns):
4 #   Column                Non-Null Count  Dtype
5 ---  ---
6 0   sepal length (cm)      150 non-null   float64
7 1   sepal width (cm)       150 non-null   float64
8 2   petal length (cm)      150 non-null   float64
9 3   petal width (cm)       150 non-null   float64
10 4   species                150 non-null   category
11 dtypes: category(1), float64(4)
12 memory usage: 5.0 KB

```

### 3.2 Summary Statistics

The output of `df.describe()` can be converted into a LaTeX table:

Table 1: Summary statistics for the Iris dataset

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.057333	3.758000	1.199333
std	0.828066	0.435866	1.765298	0.762238
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000

### 3.3 Data Distribution by Species

We can analyze the distribution of samples by species:

```

1 # Count of samples by species
2 species_counts = df['species'].value_counts()
3 print(species_counts)

```

Listing 4: Distribution by species

This would output:

Table 2: Distribution of samples by species

Species	Count
setosa	50
versicolor	50
virginica	50

## 4 Data Manipulation Techniques

### 4.1 Selecting and Filtering Data

```

1 # Select specific columns
2 sepal_data = df[['sepal length (cm)', 'sepal width (cm)']]
3
4 # Filter data by species
5 setosa = df[df['species'] == 'setosa']
6 versicolor = df[df['species'] == 'versicolor']
7 virginica = df[df['species'] == 'virginica']
8
9 # Display the first few rows of setosa data
10 print(setosa.head())

```

Listing 5: Selecting and filtering data

## 4.2 Adding New Columns

```
1 # Create a new column for sepal area
2 df['sepal_area'] = df['sepal length (cm)'] * df['sepal width (cm)']
3
4 # Create a new column for petal area
5 df['petal_area'] = df['petal length (cm)'] * df['petal width (cm)']
6
7 # Create a column for sepal length to width ratio
8 df['sepal_ratio'] = df['sepal length (cm)'] / df['sepal width (cm)']
9
10 # Display the first few rows with new columns
11 print(df.head())
```

Listing 6: Creating new features

## 4.3 Grouping and Aggregation

```
1 # Group by species and calculate mean of each feature
2 species_means = df.groupby('species').mean()
3 print(species_means)
4
5 # Multiple aggregations
6 aggs = df.groupby('species').agg({
7     'sepal length (cm)': ['min', 'max', 'mean', 'std'],
8     'petal length (cm)': ['min', 'max', 'mean', 'std']
9 })
10 print(aggs)
```

Listing 7: Grouping and aggregating data

The species means table would look like:

Table 3: Mean values of features by species

species	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	sepal_area	petal_area
setosa	5.006	3.428	1.462	0.246	17.161	0.361
versicolor	5.936	2.770	4.260	1.326	16.444	5.661
virginica	6.588	2.974	5.552	2.026	19.595	11.270

## 5 Statistical Analysis

### 5.1 Correlation Analysis

```
1 # Calculate correlation matrix
2 correlation_matrix = df.corr()
3 print(correlation_matrix)
```

Listing 8: Correlation analysis

The correlation matrix can be presented as:

Table 4: Correlation matrix of Iris features

	sepal length	sepal width	petal length	petal width	sepal_area	petal_area
sepal length	1.000	-0.118	0.872	0.818	0.814	0.876
sepal width	-0.118	1.000	-0.428	-0.366	0.452	-0.417
petal length	0.872	-0.428	1.000	0.963	0.506	0.983
petal width	0.818	-0.366	0.963	1.000	0.501	0.979
sepal_area	0.814	0.452	0.506	0.501	1.000	0.521
petal_area	0.876	-0.417	0.983	0.979	0.521	1.000

## 5.2 Descriptive Statistics by Group

```
1 # Calculate descriptive statistics by species
2 stats_by_species = df.groupby('species').describe()
3 print(stats_by_species)
```

Listing 9: Descriptive statistics by species

## 6 Data Visualization with pandas

Here are some common visualization techniques that can be used with pandas and subsequently included in LaTeX documents:

### 6.1 Histograms

```
1 import matplotlib.pyplot as plt
2
3 # Create histograms for sepal length by species
4 plt.figure(figsize=(10, 6))
5 for species in iris.species.unique():
6     subset = iris[iris.species == species]
7     plt.hist(subset['sepal length (cm)'], alpha=0.5, label=species)
8 plt.legend()
9 plt.xlabel('Sepal Length (cm)')
10 plt.ylabel('Frequency')
11 plt.title('Distribution of Sepal Length by Species')
12 plt.savefig('sepal_length_hist.png', dpi=300, bbox_inches='tight')
13 plt.close()
```

Listing 10: Creating histograms

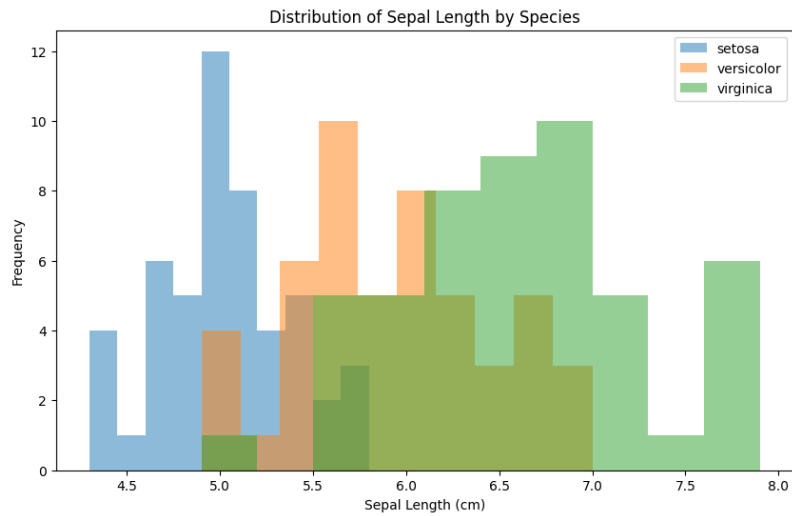


Figure 2: Distribution of sepal length by species

## 6.2 Box Plots

```

1 # Create box plots for all features by species
2 plt.figure(figsize=(12, 8))
3 df.boxplot(by='species', figsize=(12, 8))
4 plt.savefig('boxplots.png', dpi=300, bbox_inches='tight')
5 plt.close()

```

Listing 11: Creating box plots



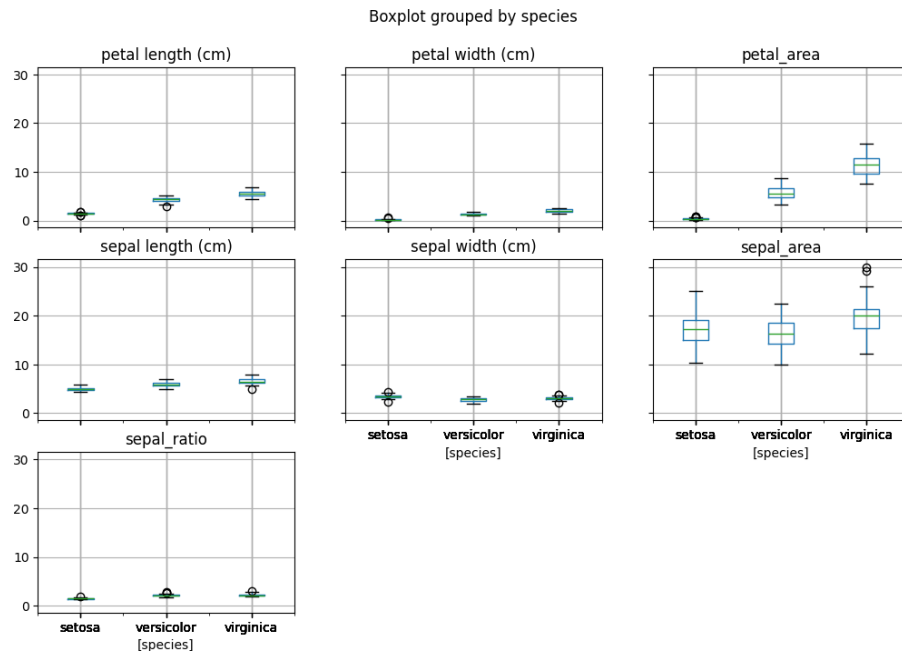


Figure 3: Box plots of Iris features by species

### 6.3 Scatter Plots

```

1 # Create scatter plot of sepal length vs width
2 plt.figure(figsize=(8, 6))
3 for species, color in zip(iris.species.unique(), ['blue', 'green', 'red']):
4     subset = iris[iris.species == species]
5     plt.scatter(subset['sepal length (cm)'],
6                 subset['sepal width (cm)'],
7                 c=color,
8                 label=species,
9                 alpha=0.7)
10 plt.xlabel('Sepal Length (cm)')
11 plt.ylabel('Sepal Width (cm)')
12 plt.title('Sepal Length vs Width by Species')
13 plt.legend()
14 plt.grid(True, linestyle='--', alpha=0.7)
15 plt.savefig('sepal_scatter.png', dpi=300, bbox_inches='tight')
16 plt.close()

```

Listing 12: Creating scatter plots

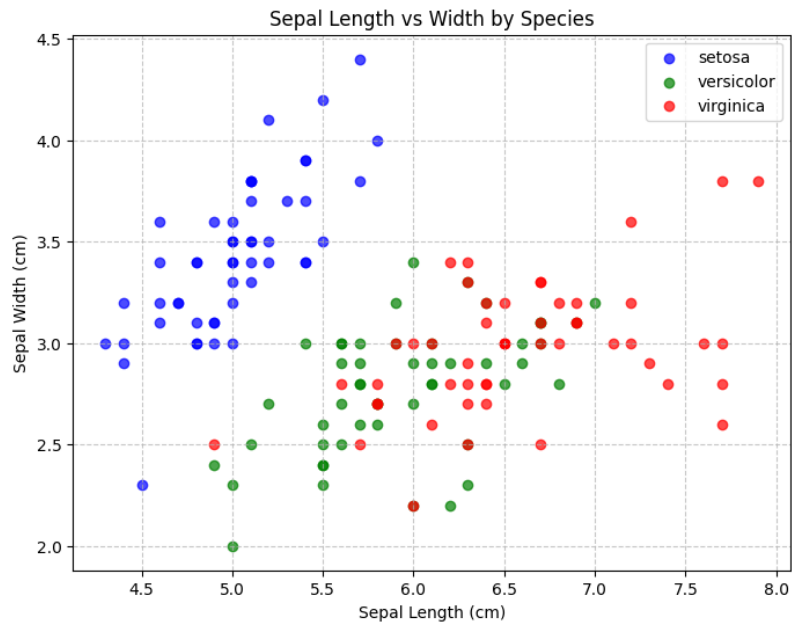


Figure 4: Scatter plot of sepal length vs width by species

## 7 Advanced pandas Techniques

### 7.1 Pivot Tables

```

1 # Create a pivot table
2 pivot = pd.pivot_table(df,
3                         values=['sepal length (cm)', 'petal length (cm)'],
4                         index='species',
5                         aggfunc=['mean', 'std'])
6 print(pivot)

```

Listing 13: Creating pivot tables

The pivot table would look like:

Table 5: Pivot table of sepal and petal length by species

species	mean		std	
	sepal length (cm)	petal length (cm)	sepal length (cm)	petal length (cm)
setosa	5.006	1.462	0.352	0.174
versicolor	5.936	4.260	0.516	0.470
virginica	6.588	5.552	0.636	0.552

## 7.2 MultiIndex and Advanced Indexing

```
1 # Create a MultiIndex DataFrame
2 arrays = [
3     ['setosa', 'setosa', 'versicolor', 'versicolor', 'virginica', 'virginica'],
4     ['sepal', 'petal', 'sepal', 'petal', 'sepal', 'petal']
5 ]
6 index = pd.MultiIndex.from_arrays(arrays, names=('species', 'part'))
7 columns = ['length', 'width']
8 data = [
9     [5.006, 3.428], [1.462, 0.246],
10    [5.936, 2.770], [4.260, 1.326],
11    [6.588, 2.974], [5.552, 2.026]
12 ]
13 multi_df = pd.DataFrame(data, index=index, columns=columns)
14 print(multi_df)
```

Listing 14: Using MultiIndex

This produces a MultiIndex DataFrame:

Table 6: MultiIndex representation of Iris measurements

species	part	length	width
setosa	sepal	5.006	3.428
	petal	1.462	0.246
versicolor	sepal	5.936	2.770
	petal	4.260	1.326
virginica	sepal	6.588	2.974
	petal	5.552	2.026

## 7.3 Applying Functions with apply()

```
1 # Define a custom function to calculate the ratio
2 def calculate_ratio(row):
3     return row['petal length (cm)'] / row['sepal length (cm)']
4
5 # Apply the function to each row
6 df['petal_sepal_ratio'] = df.apply(calculate_ratio, axis=1)
7
8 # Calculate summary statistics by species
9 ratio_by_species = df.groupby('species')['petal_sepal_ratio'].describe()
10 print(ratio_by_species)
```

Listing 15: Using apply() for custom functions

The results would look like:

Table 7: Summary statistics of petal-to-sepal length ratio by species

species	count	mean	std	min	25%	50%	75%	max
setosa	50	0.293	0.039	0.209	0.269	0.291	0.315	0.378
versicolor	50	0.717	0.050	0.581	0.686	0.720	0.750	0.815
virginica	50	0.842	0.049	0.738	0.810	0.845	0.877	0.958

## 8 Exporting Data for LaTeX

### 8.1 Converting DataFrames to LaTeX Tables

pandas provides built-in functionality to convert DataFrames to LaTeX tables:

```

1 # Convert DataFrame to LaTeX table
2 latex_table = df.head().to_latex(index=False)
3 print(latex_table)
4
5 # Save to file
6 with open('iris_table.tex', 'w') as f:
7     f.write(latex_table)

```

Listing 16: Converting DataFrames to LaTeX

### 8.2 Customizing LaTeX Output

```

1 # Customize LaTeX output
2 latex_table = df.head().to_latex(
3     index=False,
4     float_format="%.2f",
5     label="tab:iris_sample",
6     caption="Sample of the Iris dataset",
7     position="htbp",
8     column_format="lrrrrl"
9 )
10 print(latex_table)

```

Listing 17: Customizing LaTeX output

## 9 Working with Missing Data

Although the Iris dataset doesn't contain missing values, handling missing data is a common task in data analysis:

```

1 # Create a copy with some missing values
2 df_missing = df.copy()
3 df_missing.loc[0:5, 'sepal length (cm)'] = np.nan
4 df_missing.loc[10:15, 'petal width (cm)'] = np.nan
5
6 # Check for missing values
7 print(df_missing.isnull().sum())

```

```

8
9 # Drop rows with missing values
10 df_no_missing = df_missing.dropna()
11
12 # Fill missing values with mean
13 df_filled = df_missing.fillna(df.mean())
14
15 # Fill missing values with different methods per column
16 df_filled2 = df_missing.fillna({
17     'sepal length (cm)': df['sepal length (cm)'].mean(),
18     'petal width (cm)': df['petal width (cm)'].median()
19 })

```

Listing 18: Handling missing data

## 10 Merging and Joining DataFrames

```

1 # Create two DataFrames
2 measurements = df[['sepal length (cm)', 'sepal width (cm)',
3                   'petal length (cm)', 'petal width (cm)']]
4 measurements.index.name = 'id'
5 measurements = measurements.reset_index()
6
7 species_df = df[['species']].copy()
8 species_df.index.name = 'id'
9 species_df = species_df.reset_index()
10
11 # Merge DataFrames
12 merged_df = pd.merge(measurements, species_df, on='id')
13 print(merged_df.head())
14
15 # Different types of joins
16 left_join = pd.merge(measurements, species_df, on='id', how='left')
17 right_join = pd.merge(measurements, species_df, on='id', how='right')
18 outer_join = pd.merge(measurements, species_df, on='id', how='outer')

```

Listing 19: Merging DataFrames

## 11 Advanced Filtering and Selection

```

1 # Filter with complex conditions
2 large_sepals = df[(df['sepal length (cm)'] > 6.0) &
3                  (df['sepal width (cm)'] > 3.0)]
4 print(large_sepals.head())
5
6 # Use query method
7 large_sepals_query = df.query('sepal length (cm) > 6.0 and sepal width (cm) >
8                               3.0')
9 print(large_sepals_query.head())
10
11 # Select with loc and iloc
12 # loc uses labels
13 subset_loc = df.loc[10:15, ['sepal length (cm)', 'species']]

```

```

13 print(subset_loc)
14
15 # iloc uses integer positions
16 subset_iloc = df.iloc[10:15, [0, 4]]
17 print(subset_iloc)

```

Listing 20: Advanced filtering

## 12 Pandas with Iris Dataset: Case Studies

### 12.1 Case Study 1: Species Classification

```

1 # Create a model based on simple thresholds
2 def classify_iris(row):
3     if row['petal length (cm)'] < 2.5:
4         return 'setosa'
5     elif row['petal width (cm)'] < 1.8:
6         return 'versicolor'
7     else:
8         return 'virginica'
9
10 # Apply the classification function
11 df['predicted_species'] = df.apply(classify_iris, axis=1)
12
13 # Calculate accuracy
14 accuracy = (df['species'] == df['predicted_species']).mean()
15 print(f"Classification accuracy: {accuracy:.2f}")
16
17 # Create a confusion matrix
18 conf_matrix = pd.crosstab(df['species'], df['predicted_species'],
19                           rownames=['Actual'], colnames=['Predicted'])
20 print(conf_matrix)

```

Listing 21: Species classification with basic metrics

### 12.2 Case Study 2: Feature Engineering

```

1 # Create new features
2 df['sepal_area'] = df['sepal length (cm)'] * df['sepal width (cm)']
3 df['petal_area'] = df['petal length (cm)'] * df['petal width (cm)']
4 df['sepal_perimeter'] = 2 * (df['sepal length (cm)'] + df['sepal width (cm)'])
5 df['petal_perimeter'] = 2 * (df['petal length (cm)'] + df['petal width (cm)'])
6 df['ratio_areas'] = df['petal_area'] / df['sepal_area']
7
8 # Check which features best separate the species
9 for feature in ['sepal_area', 'petal_area', 'sepal_perimeter',
10                'petal_perimeter', 'ratio_areas']:
11     print(f"\nFeature: {feature}")
12     print(df.groupby('species')[feature].describe()[['mean', 'std']])

```

Listing 22: Feature engineering

## 13 Conclusion

This guide has covered a comprehensive set of pandas techniques for working with the Iris dataset, from basic data manipulation to advanced analytics. The examples provided showcase how to analyze, visualize, and export data in ways that are compatible with LaTeX documents. By integrating pandas with LaTeX, researchers and data scientists can create highly professional reports and publications that combine statistical analysis with proper documentation.

## A Complete Code Example

Here's a complete example that combines many of the techniques discussed in this guide:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn import datasets
6
7 # Load the Iris dataset
8 iris = datasets.load_iris()
9 df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
10 df['species'] = pd.Categorical.from_codes(iris.target, iris.target_names)
11
12 # Basic exploratory analysis
13 print("Dataset shape:", df.shape)
14 print("\nFirst 5 rows:")
15 print(df.head())
16 print("\nDataset info:")
17 print(df.info())
18 print("\nSummary statistics:")
19 print(df.describe())
20
21 # Feature engineering
22 df['sepal_area'] = df['sepal length (cm)'] * df['sepal width (cm)']
23 df['petal_area'] = df['petal length (cm)'] * df['petal width (cm)']
24 df['petal_sepal_ratio'] = df['petal length (cm)'] / df['sepal length (cm)']
25
26 # Group statistics
27 print("\nMean values by species:")
28 species_means = df.groupby('species').mean()
29 print(species_means)
30
31 # Visualization
32 plt.figure(figsize=(10, 8))
33 sns.pairplot(df, hue='species')
34 plt.savefig('iris_pairplot.png', dpi=300, bbox_inches='tight')
35 plt.close()
36
37 # Create a scatter plot with sepal vs petal area
38 plt.figure(figsize=(8, 6))
39 for species, color in zip(df['species'].unique(), ['blue', 'green', 'red']):
40     subset = df[df['species'] == species]
41     plt.scatter(subset['sepal_area'], subset['petal_area'],
42                 c=color, label=species, alpha=0.7)
43 plt.xlabel('Sepal Area (cm)')
```

```

44 plt.ylabel('Petal Area (cm )')
45 plt.title('Sepal Area vs Petal Area by Species')
46 plt.legend()
47 plt.grid(True, linestyle='--', alpha=0.7)
48 plt.savefig('area_scatter.png', dpi=300, bbox_inches='tight')
49 plt.close()
50
51 # Export results to LaTeX
52 with open('iris_analysis.tex', 'w') as f:
53     f.write("\\section{Iris Dataset Analysis Results}\\n\\n")
54
55     # Write summary statistics
56     f.write("\\subsection{Summary Statistics}\\n")
57     f.write(df.describe().to_latex(float_format="%.3f"))
58
59     # Write species means
60     f.write("\\n\\subsection{Mean Values by Species}\\n")
61     f.write(species_means.to_latex(float_format="%.3f"))
62
63     # Include figures
64     f.write("\\n\\subsection{Visualizations}\\n")
65     f.write("\\begin{figure}[htbp]\\n")
66     f.write("    \\centering\\n")
67     f.write("    \\includegraphics[width=0.8\\textwidth]{iris_pairplot.png}\\n")
68     f.write("    \\caption{Pairplot of Iris features by species}\\n")
69     f.write("    \\label{fig:pairplot}\\n")
70     f.write("\\end{figure}\\n\\n")
71
72     f.write("\\begin{figure}[htbp]\\n")
73     f.write("    \\centering\\n")
74     f.write("    \\includegraphics[width=0.7\\textwidth]{area_scatter.png}\\n")
75     f.write("    \\caption{Scatter plot of sepal area vs petal area by species}\\n")
76     f.write("    \\label{fig:area_scatter}\\n")
77     f.write("\\end{figure}\\n")
78
79 print("Analysis complete. Results exported to LaTeX files.")

```

Listing 23: Complete analysis workflow