

Data Analysis Pipeline Documentation

Pramit, Shraddha, Safna

December 17, 2024

Contents

1	Introduction	3
2	Setup	3
3	Function Definitions	4
3.1	Reading CSV Files	4
3.2	User Input for Wavelength Range	5
3.3	Dataset Preparation	6
3.4	Statistical calculation and graph plot for the dataset	8
3.5	For visualizing the wavelength vs absorption data	10
3.6	Pattern detection using monotonicity and correlations	11
3.7	Data smoothing	12
3.8	T-test relation between balnks and samples	13
3.9	Principal Component Analysis (PCA)	15
3.10	Get valid segment count	16
4	Main Execution Flow	17

1 Introduction

This document provides a detailed explanation of a Python-based data analysis pipeline. It includes step-by-step descriptions of the functions used, their arguments, and outputs.

2 Setup

Before using the pipeline, ensure you have the required libraries installed:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from scipy.stats import spearmanr, ttest_ind
6 from scipy.signal import savgol_filter
7 from sklearn.decomposition import PCA
```

Explanation:

- **Line 1:** `import numpy as np` The NumPy library is imported and aliased as `np`. It provides support for numerical computations and efficient handling of arrays and matrices, such as `np.array()` and `np.mean()`.
- **Line 2:** `import pandas as pd` The Pandas library is imported and aliased as `pd`. It is used for handling structured data, particularly for creating and manipulating `DataFrame` objects and performing operations like filtering, aggregation, and transformation.
- **Line 3:** `import matplotlib.pyplot as plt` The `pyplot` module of the `Matplotlib` library is imported and aliased as `plt`. This module is used for data visualization, such as creating scatter plots, line graphs, and bar charts.
- **Line 4:** `import seaborn as sns` The Seaborn library is imported as `sns`. It extends `Matplotlib` to simplify the creation of statistically informative visualizations like heatmaps, pair plots, and violin plots.
- **Line 5:** `from scipy.stats import spearmanr, ttest_ind` Specific functions from `scipy.stats` are imported:
 - `spearmanr`: Calculates the Spearman rank correlation, useful for measuring the monotonic relationship between two variables.
 - `ttest_ind`: Performs an independent two-sample t-test to compare the means of two datasets.
- **Line 6:** `from scipy.signal import savgol_filter` The Savitzky-Golay filter is imported from `scipy.signal`. This filter is used for smoothing data by fitting successive polynomial curves, preserving important trends while reducing noise.
- **Line 7:** `from sklearn.decomposition import PCA` The PCA (Principal Component Analysis) module from `Scikit-learn` is imported. PCA is a dimensionality reduction technique that transforms high-dimensional data into a lower-dimensional space while preserving as much variance as possible.

3 Function Definitions

3.1 Reading CSV Files

```
1 def processCsv(filepath):
2     """
3     Reads a CSV file and returns a pandas DataFrame.
4
5     Args:
6         filepath (str): The path to the CSV file.
7
8     Returns:
9         DataFrame: The loaded dataset as a pandas DataFrame.
10    """
11    try:
12        df = pd.read_csv(filepath)
13    except FileNotFoundError:
14        print(f"File not found: {filepath}")
15        return None
16    except Exception as e:
17        print(f"Error reading CSV file: {filepath}")
18        print(e)
19        return None
20    return df
```

Explanation:

- **Function Definition:** `def processCsv(filepath):` This defines a function named `processCsv` that takes a single parameter, `filepath`, which is a string representing the path to the CSV file.
- **Docstring:** The docstring provides a description of the function's purpose, arguments, and return value.
 - **Args:** `filepath (str)`: Specifies the path to the CSV file to be read.
 - **Returns:** `DataFrame`: A pandas DataFrame containing the data from the CSV file.
- **Try Block:** The `try` block attempts to execute the code to read the CSV file:
 - `df = pd.read_csv(filepath)`: This uses `pandas.read_csv()` to load the CSV file into a pandas DataFrame.
- **Exception Handling:** The `except` blocks handle errors that may occur while reading the file:
 - `FileNotFoundError`: If the file at the specified path does not exist, an error message is printed: `"File not found: {filepath}"`, and `None` is returned.
 - `Exception as e`: Captures any other exceptions, prints an error message, and returns `None`.
- **Return Statement:** If no exception occurs, the loaded DataFrame `df` is returned as the output of the function.

3.2 User Input for Wavelength Range

```
1 def getWavelengthRange():
2     """
3     Prompt the user to input the desired wavelength range for
4     segmentation.
5
6     Returns:
7         tuple: A tuple containing the lower and upper bounds of
8         the wavelength range.
9     """
10    while True:
11        try:
12            lower_bound = float(input("Enter the lower bound of
13            the wavelength range: "))
14            upper_bound = float(input("Enter the upper bound of
15            the wavelength range: "))
16            print("Lower bound:", lower_bound)
17            print("Upper bound:", upper_bound)
18            if lower_bound >= upper_bound:
19                print("Lower bound must be less than the upper
20                bound. Please try again.")
21            else:
22                return lower_bound, upper_bound
23        except ValueError:
24            print("Invalid input. Please enter numeric values.")
```

Explanation:

- **Function Definition:** `def getWavelengthRange():` This defines a function named `getWavelengthRange` that prompts the user for input and returns a tuple representing the lower and upper bounds of a wavelength range.
- **Docstring:** The docstring provides a brief overview of the function's purpose and its return value.
 - **Returns:** tuple: A tuple containing two floating-point values (`lower_bound`, `upper_bound`) that represent the wavelength range.
- **While Loop:** `while True:` This loop ensures the user is continuously prompted until valid input is provided. It keeps the program running until the user enters acceptable values.
- **Try Block:** The `try` block handles user input and potential errors:
 - `lower_bound = float(input(...))`: Prompts the user to input the lower bound of the wavelength range, converting the input to a float.
 - `upper_bound = float(input(...))`: Prompts the user to input the upper bound of the wavelength range, converting the input to a float.
- **Validation Check:** `if lower_bound >= upper_bound:` Ensures that the lower bound is less than the upper bound. If not, an error message is printed:

- "Lower bound must be less than the upper bound. Please try again."
- **Successful Input:** return lower_bound, upper_bound If the inputs are valid and the lower bound is less than the upper bound, the function returns a tuple containing the two values.
- **Exception Handling:** except ValueError: If the user inputs a non-numeric value, a ValueError is raised. An error message is displayed:
 - "Invalid input. Please enter numeric values."

3.3 Dataset Preparation

```

1 def prepareDataset(df, filepath):
2     """
3     Cleans and preprocesses the dataset by renaming columns,
4     filtering rows, and converting data types.
5
6     Args:
7         df (DataFrame): The raw dataset.
8         filepath (str): The path to the CSV file (used for error
9         logging).
10
11     Returns:
12         DataFrame: A cleaned and preprocessed dataset.
13     """
14     print("Data cleaning....")
15     column_names = list(df.iloc[0]) # Extract column names from
16     the first row
17     column_names[0] = "Wavelength (nm)" # Rename the first
18     column for clarity
19     filtered_data = df[1:].copy() # Exclude the first row used
20     for column names
21     filtered_data.columns = column_names # Assign column names
22
23     # Convert "Wavelength (nm)" to numeric
24     try:
25         filtered_data["Wavelength (nm)"] = pd.to_numeric(
26             filtered_data["Wavelength (nm)"], errors='coerce')
27     except Exception as e:
28         print(f"Error converting 'Wavelength (nm)' to numeric in
29             file: {filepath}")
30         print(e)
31         return None
32
33     # Drop rows with invalid "Wavelength (nm)" values
34     filtered_data = filtered_data.dropna(subset=["Wavelength (nm)
35         "])
36
37     # Get wavelength range from the user

```

```

30     print("Please specify the wavelength range for segmentation."
31           )
32     lower_bound, upper_bound = getWavelengthRange()
33
34     # Filter data within the user-specified range
35     filtered_data = filtered_data[(filtered_data["Wavelength (nm)"]
36                                     >= lower_bound) &
37                                     (filtered_data["Wavelength (nm)"]
38                                     <= upper_bound)]
39
40     # Convert absorption columns to numeric
41     valid_columns = ["Wavelength (nm)"] # Keep track of valid
42     columns
43     for col in filtered_data.columns[1:]:
44         try:
45             filtered_data[col] = pd.to_numeric(filtered_data[col], errors='coerce')
46             valid_columns.append(col)
47         except Exception as e:
48             print(f"Skipping invalid column '{col}' in file: {filepath}")
49             print(e)
50
51     # Keep only valid columns
52     filtered_data = filtered_data[valid_columns]
53     print("Data is prepared now!")
54     return filtered_data

```

Explanation:

- **Function Definition:** `def prepareDataset(df, filepath):` This function cleans and preprocesses a dataset provided as a `pandas DataFrame`, using the file path for logging purposes in case of errors.
- **Docstring:** The docstring explains the function's purpose, arguments, and return value.
 - **Args:** `df (DataFrame)`: The raw dataset to be cleaned and processed. `filepath (str)`: The path to the CSV file, used for error messages.
 - **Returns:** `DataFrame`: A cleaned and preprocessed dataset.
- **Column Renaming and Filtering:**
 - `column_names = list(df.iloc[0]):` Extracts column names from the first row of the dataset.
 - `column_names[0] = "Wavelength (nm)":` Renames the first column for clarity.
 - `filtered_data = df[1:].copy():` Excludes the first row, which contains the column names, and creates a copy of the remaining data.
 - `filtered_data.columns = column_names:` Assigns the extracted column names to the dataset.

- **Numeric Conversion of Wavelength Column:**

- `filtered_data["Wavelength (nm)"] = pd.to_numeric(...)`: Converts the "Wavelength (nm)" column to numeric, coercing invalid values to NaN.
- If an error occurs during conversion, the exception is caught, and an error message is printed with the file path.

- **Dropping Invalid Rows:**

- `filtered_data = filtered_data.dropna(subset=["Wavelength (nm)"])`: Removes rows where "Wavelength (nm)" contains NaN.

- **Wavelength Range Filtering:**

- The `getWavelengthRange()` function is called to prompt the user for the lower and upper bounds of the wavelength range.
- `filtered_data = filtered_data[(filtered_data["Wavelength (nm)"] >= lower_bound) & (filtered_data["Wavelength (nm)"] <= upper_bound)]`: Filters the dataset to include only rows within the specified range.

- **Numeric Conversion of Absorption Columns:**

- A list `valid_columns` is initialized with "Wavelength (nm)".
- For each column other than "Wavelength (nm)", the function attempts to convert it to numeric using `pd.to_numeric(...)`.
- If a column cannot be converted, an error message is printed, and the column is skipped.

- **Keeping Only Valid Columns:**

- `filtered_data = filtered_data[valid_columns]`: Keeps only the valid columns in the final dataset.

- **Return Statement:**

- If all operations are successful, the cleaned and preprocessed dataset is returned.

- **Output Messages:**

- Messages like "Data cleaning..." and "Data is prepared now!" provide feedback to the user during processing.

3.4 Statistical calculation and graph plot for the dataset

```
1 def statForDataset(df):
2     """
3     Computes and displays statistical summaries for the dataset.
4
5     Args:
6         df (DataFrame): The dataset to analyze.
```



```

7      """
8      print("Dataset Statistics:")
9      # Compute basic statistics for all columns except the
        wavelength
10     data_no_wavelength = df.iloc[:, 1:] # Exclude the first
        column (wavelength)
11
12     # Compute statistics
13     statistics = {
14         "Mean": data_no_wavelength.mean(),
15         "Median": data_no_wavelength.median(),
16         "Standard Deviation": data_no_wavelength.std(),
17         "Min": data_no_wavelength.min(),
18         "Max": data_no_wavelength.max(),
19         "Range": data_no_wavelength.max() - data_no_wavelength.min(),
20     }
21
22     # Print each statistic
23     for stat_name, stat_values in statistics.items():
24         print(f"\n{stat_name}:")
25         print(stat_values.to_string(index=True))
26
27     # Plotting
28     plt.figure(figsize=(15, 8)) # Increased figure size
29
30     # Create subplots for each statistic
31     for i, stat_name in enumerate(statistics, 1):
32         plt.subplot(3, 3, i)
33         statistics[stat_name].plot(kind='bar',
34                                     color='skyblue',
35                                     edgecolor='black')
36
37         plt.title(f'{stat_name} Across Columns', fontsize=12)
38         plt.xticks(rotation=45, ha='right', fontsize=8)
39         plt.ylabel(stat_name, fontsize=10)
40         plt.grid(axis='y', linestyle='--', alpha=0.7)
41
42     plt.suptitle('Statistical Measures Comparison', fontsize=16)
43     plt.tight_layout()
44     plt.show()

```

Explanation:

- **Function Definition and Purpose:**

- Function `statForDataset(df)` is designed to perform comprehensive statistical analysis on a pandas DataFrame
- Takes a DataFrame as input and excludes the first column (assumed to be wavelength)

- **Key Code Analysis:**

- `df.iloc[:, 1:]`: Slices DataFrame, removing first column

- Statistical computations use pandas built-in methods:
 - * `.mean()`: Calculates column-wise average
 - * `.median()`: Finds middle value of each column
 - * `.std()`: Computes standard deviation
 - * `.min()/max()`: Determines minimum/maximum values
- **Statistical Dictionary Creation:**
 - Creates `statistics` dictionary with multiple statistical measures
 - Enables flexible, comprehensive data exploration
 - Computes range using `max() - min()` subtraction
- **Visualization Approach:**
 - Uses matplotlib for creating subplots
 - `plt.subplot(3, 3, i)`: Creates 3x3 grid of statistical visualizations
 - Bar plot with customized aesthetics:
 - * Color: skyblue
 - * Edge color: black
 - * Rotated x-axis labels for readability

3.5 For visualizing the wavelength vs absorption data

```

1 def visualizeTrends(df):
2     """
3     Plots the absorption trends for each column in the dataset.
4
5     Args:
6         df (DataFrame): The dataset containing wavelength and
7             absorption data.
8     """
9     plt.figure(figsize=(15, 8))
10    for col in df.columns[1:]:
11        plt.plot(df["Wavelength (nm)"], df[col], label=col)
12    plt.title("Absorption Trends")
13    plt.xlabel("Wavelength (nm)")
14    plt.ylabel("Absorption")
15    plt.legend()
16    plt.grid()
17    plt.show()

```

Explanation:

- **Function Signature:** `visualizeTrends(df)` A specialized data visualization function designed to plot absorption trends across multiple spectroscopic measurements.
- **Function Objectives:**

- Generate comprehensive absorption trend visualization
 - Display multiple dataset columns simultaneously
 - Provide comparative spectral analysis
- **Visualization Methodology:**
 - **Matplotlib Configuration:**
 - * `plt.figure(figsize=(15, 8))`: Large visualization canvas
 - * Enables detailed, high-resolution plotting
 - **Plotting Strategy:**
 - * Iterates through all columns except first (wavelength)
 - * `df.columns[1:]`: Selects absorption data columns
 - * Uses `plt.plot()` for line representation
 - **Graphical Elements:**
 - X-axis: Wavelength (nm)
 - Y-axis: Absorption intensity
 - Automatic legend generation
 - Grid lines for enhanced readability

3.6 Pattern detection using monotonicity and correlations

```

1 def patternDetection(df):
2     """
3     Analyzes patterns using monotonicity (Spearman correlation)
4     and visualizes inter-column correlations.
5
6     Args:
7         df (DataFrame): The dataset for analysis.
8     """
9     print("Monotonicity and Correlation Analysis:")
10
11     # Compute monotonicity using Spearman's rank correlation
12     monotonicity_results = {}
13     for col in df.columns[1:]:
14         monotonicity, _ = spearmanr(df["Wavelength (nm)"], df[col])
15         monotonicity_results[col] = monotonicity
16         print(f"Monotonicity (Spearman's rho) for {col}: {monotonicity:.2f}")
17
18     # Plot correlation heatmap
19     corr_df = df.iloc[:, 1:] # drop the 'Wavelength (nm)'
20     correlations = corr_df.corr(method='spearman')
21     plt.figure(figsize=(15, 8))
22     sns.heatmap(correlations, annot=True, cmap="coolwarm", vmin=-1, vmax=1)

```

```

22 plt.title("Correlation Heatmap")
23 plt.show()

```

Explanation

- **Function Signature:** `patternDetection(df)` A sophisticated statistical analysis function designed to detect and visualize complex data patterns using advanced correlation techniques.
- **Analytical Objectives:**
 - Compute monotonicity using Spearman's rank correlation
 - Visualize inter-column correlational relationships
 - Provide comprehensive pattern detection insights
- **Monotonicity Analysis:**
 - **Computational Approach:**
 - * Uses `spearmanr()` for non-linear correlation detection
 - * Compares wavelength against each absorption column
 - * Captures monotonic trends independent of linear relationships
 - **Correlation Interpretation:**
 - * Range: -1 to 1
 - * 0: No monotonic relationship
 - * Close to ± 1 : Strong monotonic trend
- **Visualization Techniques:**
 - **Correlation Heatmap:**
 - * Uses seaborn's `heatmap()` for visualization
 - * Color scheme: "coolwarm" (-1 to 1 range)
 - * Annotated correlation coefficients
 - **Graphical Configuration:**
 - * `figsize=(15, 8)`: Large, detailed visualization
 - * Comprehensive inter-column correlation representation

3.7 Data smoothing

```

1 def dataSmoothing(df):
2     """
3     Applies Savitzky-Golay filtering to smooth absorption data.
4
5     Args:
6         df (DataFrame): The dataset to smooth.
7
8     Returns:
9         DataFrame: The smoothed dataset.

```

```

10     """
11     smoothed_data = df.copy()
12     # Exclude the first column from smoothing
13     for col in smoothed_data.columns[1:]:
14         smoothed_data[col] = savgol_filter(smoothed_data[col],
15                                             window_length=11, polyorder=3)
16     if len(smoothed_data.columns) > 1:
17         second_column_name = smoothed_data.columns[1]
18         filtered_data = smoothed_data.drop(columns=[
19             second_column_name])
20         print(f"Dropped second column: {second_column_name}")
21     visualizeTrends(filtered_data)
22     return smoothed_data

```

Explanation

- **Function Signature:** `dataSmoothing(df)` A sophisticated signal processing function designed to apply advanced noise reduction techniques on spectroscopic datasets.
- **Smoothing Methodology:**
 - **Algorithm:** Savitzky-Golay Filter
 - * Preserves higher-order moments
 - * Maintains signal characteristics
 - * Reduces random noise
 - **Filter Parameters:**
 - * `window_length=11`: Sliding window size
 - * `polyorder=3`: Polynomial approximation order
 - * Balances noise reduction and signal preservation
- **Implementation Strategy:**
 - Creates deep copy of input DataFrame
 - Applies smoothing to all columns except wavelength
 - Uses `savgol_filter()` for computational efficiency
- **Post-Processing Actions:**
 - Conditional column dropping mechanism
 - Calls `visualizeTrends()` for result visualization
 - Returns smoothed dataset for further analysis

3.8 T-test relation between balnks and samples

```

1 def testRelationBetweenBlanksAndSamples(df):
2     """
3     Performs t-tests to compare blanks and samples for
4     statistical differences.

```

```

5     Args:
6         df (DataFrame): The dataset containing blanks and sample
          data.
7     """
8
9     # Explicitly identify blank and sample columns based on known
        naming patterns
10    blank_columns = [col for col in df.columns if col.startswith(
        "blank")]
11    sample_columns = [col for col in df.columns if col.endswith("
        mg")]
12
13    if not blank_columns or not sample_columns:
14        print("No blank or sample columns found in the dataset.")
15        return
16
17    for blank_col in blank_columns:
18        for sample_col in sample_columns:
19            try:
20                blank_data = df[blank_col].dropna()
21                sample_data = df[sample_col].dropna()
22
23                if blank_data.empty or sample_data.empty:
24                    print(f"No data to compare between {blank_col
                } and {sample_col}.")
25                    continue
26
27                stat, p_value = ttest_ind(blank_data, sample_data
                )
28                print(f"T-test between {blank_col} and {
                sample_col}: p-value = {p_value: .4f}")
29            except Exception as e:
30                print(f"Error performing t-test between {
                blank_col} and {sample_col}: {e}")

```

Explanation

- **Function Signature:** `testRelationBetweenBlanksAndSamples(df)` A comprehensive statistical analysis function designed to compare blank and sample measurements using independent t-tests.
- **Column Identification Strategy:**
 - **Blank Column Detection:**
 - * Uses prefix-based identification
 - * `col.startswith("blank")`: Identifies blank columns
 - **Sample Column Detection:**
 - * Uses suffix-based identification
 - * `col.endswith("mg")`: Identifies sample columns
- **Statistical Testing Methodology:**

- **T-Test Configuration:**
 - * Independent two-sample t-test
 - * Compares means between blank and sample datasets
 - * Assesses statistical significance
- **Data Preparation:**
 - * Removes NaN values using `.dropna()`
 - * Ensures valid data for comparison
- **Comparative Analysis Approach:**
 - Nested loop for comprehensive cross-comparison
 - Pairwise t-test between all blank and sample columns
 - Prints p-values for statistical interpretation
- **Error Handling Mechanisms:**
 - Checks for empty datasets
 - Catches and reports exceptions during testing
 - Provides informative error messages

3.9 Principal Component Analysis (PCA)

```

1 def applyPCA(df):
2     """
3     Applies Principal Component Analysis (PCA) to reduce
4     dimensionality and visualize data in 2D.
5
6     Args:
7         df (DataFrame): The dataset for PCA.
8
9     Displays:
10        Scatter plot of the first two principal components.
11    """
12    print("Applying PCA...")
13    pca = PCA(n_components=2)
14    absorption_data = df.iloc[:, 1:].values
15    principal_components = pca.fit_transform(absorption_data)
16    pca_df = pd.DataFrame(principal_components, columns=["PC1", "PC2"])
17    print("PCA applied!And we are getting this scatter graph.")
18    plt.figure(figsize=(12, 8))
19    plt.scatter(pca_df["PC1"], pca_df["PC2"])
20    plt.title("PCA of Absorption Data")
21    plt.xlabel("Principal Component 1")
22    plt.ylabel("Principal Component 2")
23    plt.grid()
24    plt.show()

```

Explanation

- **Function Signature:** `applyPCA(df)` A dimensionality reduction function utilizing Principal Component Analysis (PCA) for spectroscopic data visualization.
- **Dimensionality Reduction Strategy:**
 - **PCA Configuration:**
 - * `n_components=2`: Reduces to 2-dimensional space
 - * Captures maximum variance in minimal dimensions
 - **Data Preparation:**
 - * `df.iloc[:, 1:]`: Excludes wavelength column
 - * Converts to numerical numpy array
- **Computational Workflow:**
 - Applies `fit_transform()` for PCA computation
 - Generates principal component DataFrame
 - Creates two-dimensional representation
- **Visualization Technique:**
 - **Scatter Plot Configuration:**
 - * X-axis: First Principal Component
 - * Y-axis: Second Principal Component
 - * Grid for enhanced readability

3.10 Get valid segment count

```
1 def getValidSegmentCount():
2     """
3     Prompt user to enter the number of segments with input
4     validation.
5
6     Returns:
7         int: Number of segments to process
8     """
9     while True:
10         try:
11             segment_count = int(input("How many segments would
12                                     you like to analyze? (1-5): "))
13
14             # Validate segment count
15             if 1 <= segment_count <= 5:
16                 return segment_count
17             else:
18                 print("Please enter a number between 1 and 5.")
19
20         except ValueError:
21             print("Invalid input. Please enter a valid integer.")
```


Explanation

- **Function Signature:** `getValidSegmentCount()` An interactive input validation function designed to collect user-specified segment count with robust error handling.
- **Input Validation Strategy:**
 - **Input Mechanism:**
 - * Uses `input()` for user interaction
 - * Prompts for segment count between 1-5
 - **Validation Techniques:**
 - * Infinite `while` loop for continuous prompting
 - * `int()` conversion for numeric validation
 - * Range check: $1 \leq \text{segment_count} \leq 5$
- **Error Handling Mechanisms:**
 - Catches `ValueError` for non-integer inputs
 - Provides user-friendly error messages
 - Ensures valid integer input
- **Return Behavior:**
 - Returns validated segment count
 - Exits loop upon successful input

4 Main Execution Flow

The main script uses the defined functions to analyze the dataset and produce visualizations:

```
1 if __name__ == "__main__":
2     filepath = '/content/set01.csv'
3     df = processCsv(filepath)
4
5     if df is not None:
6         # Get number of segments from user
7         num_segments = getValidSegmentCount()
8
9         # Lists to store segment data
10        segments = []
11        smoothed_segments = []
12
13        # Process each segment
14        for i in range(num_segments):
15            print(f"\nEnter details for segment {i+1}:")
16            current_segment = prepareDataset(df, filepath)
17
18            if current_segment is not None:
```

```

19         segments.append(current_segment)
20
21     # Visualization
22     print(f"Visualizing segment {i+1}...")
23     visualizeTrends(current_segment)
24
25     # Statistical Analysis
26     print(f"Statistical info on segment {i+1}...")
27     statForDataset(current_segment)
28
29     # Pattern Detection
30     print(f"Applying pattern detection on segment {i+1}...")
31     patternDetection(current_segment)
32
33     # Smoothing
34     print(f"Smoothing segment {i+1}...")
35     smoothed_segment = dataSmoothing(current_segment)
36     smoothed_segments.append(smoothed_segment)
37
38     # Re-apply Pattern Detection on Smoothed Data
39     print(f"Re-applying pattern detection on smoothed segment {i+1}...")
40     patternDetection(smoothed_segment)
41
42     # Test Blanks and Samples Relation
43     print(f"Testing relation between blanks and samples for segment {i+1}...")
44     testRelationBetweenBlanksAndSamples(current_segment)
45
46     # PCA
47     print(f"Applying PCA on segment {i+1}...")
48     applyPCA(current_segment)

```

Explanation

- **Main Execution Flow:** A comprehensive data analysis pipeline designed for multi-segment spectroscopic data processing.
- **Initialization Stage:**
 - **Data Loading:**
 - * Uses `processCsv()` for initial data import
 - * Validates successful CSV processing
 - **Segment Configuration:**
 - * Calls `getValidSegmentCount()` for user input
 - * Determines number of segments to analyze
- **Segment Processing Workflow:**

- **Iterative Analysis:**
 - * Processes each segment sequentially
 - * Applies multiple analytical techniques
- **Analysis Techniques:**
 - * Data Preparation
 - * Trend Visualization
 - * Statistical Analysis
 - * Pattern Detection
 - * Data Smoothing
 - * Blank-Sample Relationship Testing
 - * Principal Component Analysis
- **Data Management:**
 - Maintains separate lists for:
 - * Original segments
 - * Smoothed segments