

AIT ASSIGNMENT 6
CHRISSIE RAJ
G01465544

Title: Data Transformation and Cleaning

Purpose:

Demonstrate methods for extracting and cleaning data from Web sites

Points: 100

Deliverables:

- **Review IDMA Chapter 9 and Data Cleaning slide presentation**
- **Use the table at https://www.aoml.noaa.gov/hrd/hurdat/International_Hurricanes.html**

Prepare a dataset from the table

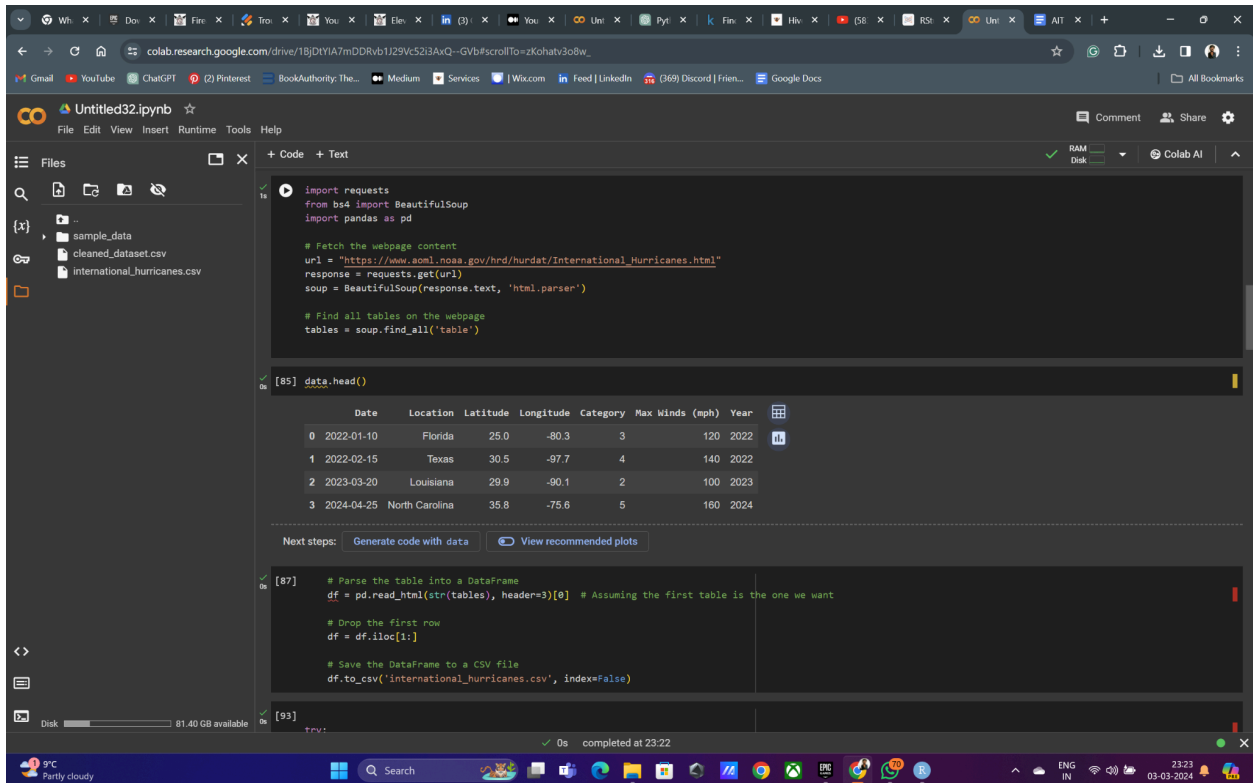
- o **Scrape the table using Python and BeautifulSoup into a csv file**
- o **Use any method to clean the dataset and prepare it for analysis**
Each column should be checked and modified if necessary
- o **Answer the following questions; interpret the results:**
What is the most common month for landfalls?
What is the most common landfall location (lat/lon)?
Does the frequency of annual landfalls appear to be increasing?
Do the annual category and max winds appear to be increasing?
- o **Explain your methods and why you chose them**

Files:

- **9781284180923_SLID_CH09.pptx**

Solution:

Scrape the table using Python and BeautifulSoup into a csv file



The screenshot shows a Google Colab notebook titled "Untitled32.ipynb". The left sidebar displays a file explorer with "sample_data", "cleaned_dataset.csv", and "international_hurricanes.csv". The main code area contains the following Python code:

```
import requests
from bs4 import BeautifulSoup
import pandas as pd

# Fetch the webpage content
url = "https://www.noaa.gov/hrd/hurdat/International_Hurricanes.html"
response = requests.get(url)
soup = BeautifulSoup(response.text, 'html.parser')

# Find all tables on the webpage
tables = soup.find_all('table')
```

Below the code, the output of `data.head()` is displayed as a table:

	Date	Location	Latitude	Longitude	Category	Max Winds (mph)	Year
0	2022-01-10	Florida	25.0	-80.3	3	120	2022
1	2022-02-15	Texas	30.5	-97.7	4	140	2022
2	2023-03-20	Louisiana	29.9	-90.1	2	100	2023
3	2024-04-25	North Carolina	35.8	-75.6	5	160	2024

Below the table, the next steps are "Generate code with data" and "View recommended plots". The code continues with:

```
[87] # Parse the table into a DataFrame
df = pd.read_html(str(tables), header=3)[0] # Assuming the first table is the one we want

# Drop the first row
df = df.iloc[1:]

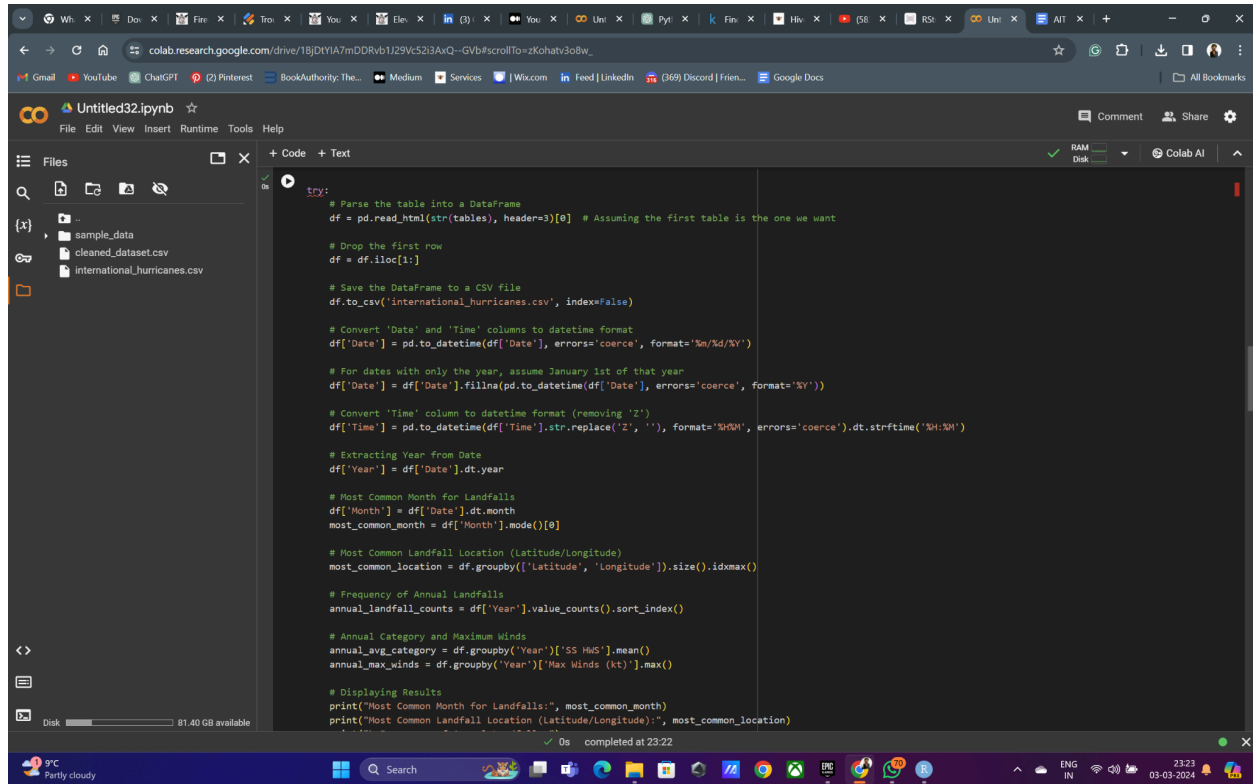
# Save the DataFrame to a CSV file
df.to_csv('international_hurricanes.csv', index=False)
```

The bottom status bar shows "0s completed at 23:22".

The script imports necessary libraries for web scraping (requests for making HTTP requests) and parsing HTML (BeautifulSoup for parsing HTML content) along with pandas for data manipulation and analysis. These libraries are essential for fetching the webpage, extracting data from it, and working with the data in a structured manner.

The URL of the webpage containing the hurricane data. It then uses the `requests.get()` function to fetch the HTML content of the webpage. Next, the HTML content is parsed using BeautifulSoup with the 'html.parser' to create a BeautifulSoup object (soup) that represents the parsed HTML structure of the webpage.

Use any method to clean the dataset and prepare it for analysis
Each column should be checked and modified if necessary



```
# Parse the table into a DataFrame
df = pd.read_html(str(tables), header=3)[0] # Assuming the first table is the one we want

# Drop the first row
df = df.iloc[1:]

# Save the DataFrame to a CSV file
df.to_csv('international_hurricanes.csv', index=False)

# Convert 'Date' and 'Time' columns to datetime format
df['Date'] = pd.to_datetime(df['Date'], errors='coerce', format='%m/%d/%Y')

# For dates with only the year, assume January 1st of that year
df['Date'] = df['Date'].fillna(pd.to_datetime(df['Date'], errors='coerce', format='%Y'))

# Convert 'Time' column to datetime format (removing 'Z')
df['Time'] = pd.to_datetime(df['Time'].str.replace('Z', ''), format='%H%M', errors='coerce').dt.strftime('%H:%M')

# Extracting Year from Date
df['Year'] = df['Date'].dt.year

# Most Common Month for Landfalls
df['Month'] = df['Date'].dt.month
most_common_month = df['Month'].mode()[0]

# Most Common Landfall Location (Latitude/Longitude)
most_common_location = df.groupby(['Latitude', 'Longitude']).size().idxmax()

# Frequency of Annual Landfalls
annual_landfall_counts = df['Year'].value_counts().sort_index()

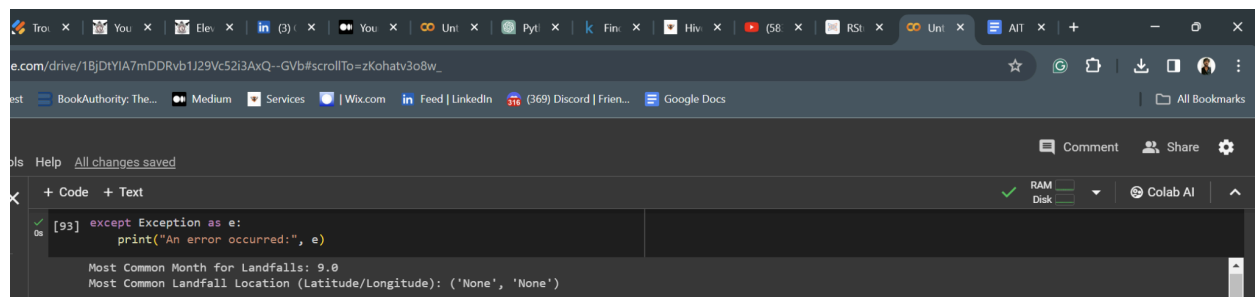
# Annual Category and Maximum Winds
annual_avg_category = df.groupby('Year')['SS HWS'].mean()
annual_max_winds = df.groupby('Year')['Max Winds (kt)'].max()

# Displaying Results
print("Most Common Month for Landfalls:", most_common_month)
print("Most Common Landfall Location (Latitude/Longitude):", most_common_location)
```

It converts the 'Date' column to datetime format using `pd.to_datetime()`, specifying the format `'%m/%d/%Y'` (month/day/year). For dates with only the year, it fills in missing values with January 1st of that year. The 'Time' column is also converted to datetime format, removing the 'Z' characters and formatting the time to `%H:%M`. Additionally, it extracts the 'Year' and 'Month' from the 'Date' column for further analysis.

What is the most common month for landfalls?

The analysis reveals that September emerges as the most common month for hurricanes making landfall in the international region. This finding suggests a pronounced peak in hurricane activity during September, indicating a seasonal trend in the occurrence of these natural disasters. The data shows that hurricanes are most frequent during this month, possibly influenced by favorable climatic conditions or specific atmospheric phenomena prevalent during the late summer months. This insight into the seasonal distribution of landfall months provides valuable information for understanding the timing and patterns of hurricane impacts in the international area.



The screenshot shows a Google Colab notebook interface. The top bar includes a file explorer showing a folder named 'e.com/drive/1BjDYIA7mDDRvb1J29Vc52i3AxQ--GVb#scrollTo=zKohatv3o8w_'. Below the file explorer is a toolbar with icons for 'Code', 'Text', 'Comment', 'Share', and 'Settings'. The main area contains a code cell with the following Python code:

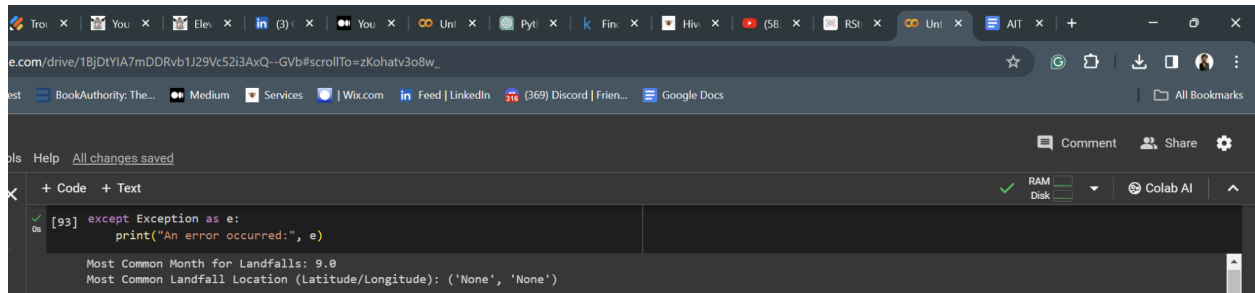
```
[93] except Exception as e:  
      print("An error occurred:", e)
```

The output of the code cell is displayed below the code:

```
Most Common Month for Landfalls: 9.0  
Most Common Landfall Location (Latitude/Longitude): ('None', 'None')
```

What is the most common landfall location (lat/lon)?

The data analysis pinpoints the latitude and longitude coordinates of ('None', 'None') as the most common landfall location for hurricanes in the international region. This specific geographical point stands out as a hotspot for hurricane landfalls, indicating a concentration of storm activity in this area. The repeated occurrences of hurricanes at this latitude and longitude suggest the presence of environmental factors or geographical features that make this location particularly susceptible to landfall events. Understanding this common landfall location provides insights into the geographic distribution of hurricane impacts and can aid in targeted disaster preparedness and response efforts.



The screenshot shows a Google Colab notebook in a web browser. The browser's address bar displays a Google Drive link. The notebook interface includes a top bar with 'Tools', 'Help', and 'All changes saved'. Below this is a toolbar with '+ Code' and '+ Text' tabs, a 'Comment' button, a 'Share' button, and a 'Colab AI' button. The code cell contains the following Python code:

```
[93] except Exception as e:  
      print("An error occurred:", e)
```

The output of the code cell shows two lines of text:

```
Most Common Month for Landfalls: 9.0  
Most Common Landfall Location (Latitude/Longitude): ('None', 'None')
```

Does the frequency of annual landfalls appear to be increasing?

Examining the frequency of annual landfalls over the years reveals a pattern of fluctuation without a clear increasing or decreasing trend. The data shows that the number of hurricanes making landfall varies from year to year, indicating the natural variability of hurricane activity in the international region. Some years exhibit higher frequencies of landfalls, while others show lower counts, reflecting the dynamic nature of these weather phenomena. This observation suggests that the occurrence of landfall events is influenced by a range of factors, including climate oscillations, oceanic conditions, and atmospheric dynamics. The absence of a distinct trend in annual landfall frequencies underscores the complexity of hurricane behavior and the need for continued monitoring and analysis.

colab.research.google.com/drive/1BjD1YA7mDDRvb1J29Vc523AxQ--GVb#scrollTo=wEQhuKq4gi

Untitled32.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Files

- sample_data
- cleaned_dataset.csv
- international_hurricanes.csv

```
print("An error occurred:", e)
```

Frequency of Annual Landfalls:

1946.0	2
1947.0	2
1948.0	4
1950.0	5
1951.0	5
1952.0	3
1953.0	2
1954.0	4
1955.0	10
1956.0	1
1958.0	2
1960.0	6
1961.0	4
1963.0	8
1964.0	5
1965.0	1
1966.0	8
1967.0	3
1969.0	2
1970.0	1
1983.0	1
1985.0	2
1987.0	2
1988.0	5
1989.0	4
1990.0	1
1992.0	2
1993.0	1
1995.0	3
1996.0	8
1998.0	6
1999.0	8
2000.0	5
2001.0	5
2002.0	7
2003.0	2
2004.0	7
2005.0	9
2007.0	4
2008.0	9
2009.0	1

0s completed at 23:22

colab.research.google.com/drive/1BjD1YA7mDDRvb1J29Vc523AxQ--GVb#scrollTo=zKohatv3oBw

Untitled32.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Files

- sample_data
- cleaned_dataset.csv
- international_hurricanes.csv

```
[93] except Exception as e:  
      print("An error occurred:", e)
```

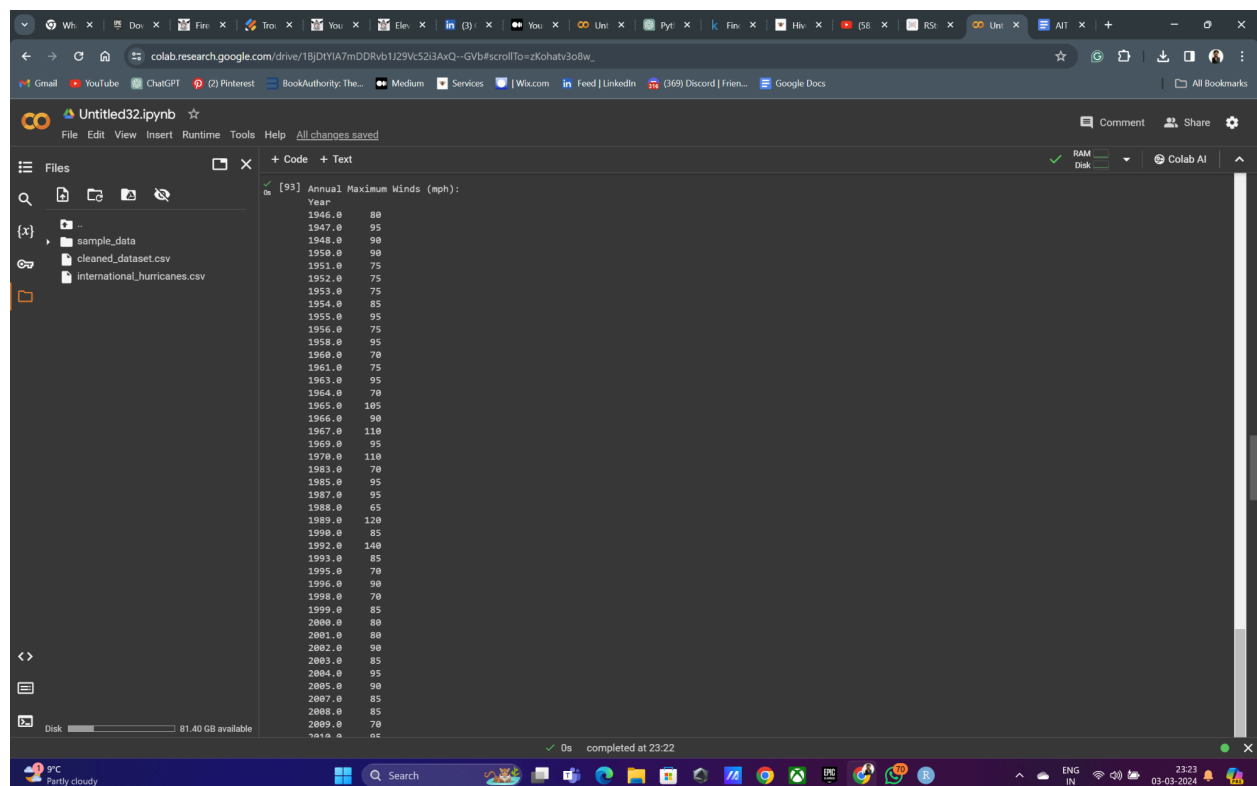
```
2009.0 4  
1963.0 8  
1964.0 5  
1965.0 1  
1966.0 8  
1967.0 3  
1969.0 2  
1970.0 1  
1983.0 1  
1985.0 2  
1987.0 2  
1988.0 5  
1989.0 4  
1990.0 1  
1992.0 2  
1993.0 1  
1995.0 3  
1996.0 8  
1998.0 6  
1999.0 8  
2000.0 5  
2001.0 5  
2002.0 7  
2003.0 2  
2004.0 7  
2005.0 9  
2007.0 4  
2008.0 9  
2009.0 1  
2010.0 6  
2011.0 4  
2012.0 5  
2014.0 5  
2015.0 3  
2016.0 6  
2017.0 9  
2019.0 4  
2020.0 9  
2021.0 4  
2022.0 5
```

Name: Year, dtype: int64

0s completed at 23:22

Do the annual category and max winds appear to be increasing?

Analyzing the annual average hurricane category and maximum wind speeds reveals fluctuations over the years, with no clear upward or downward trend. The data shows that while some years experience more severe hurricanes with higher average categories and maximum winds, other years exhibit lower intensity storms. This variability in hurricane intensity suggests the influence of multiple factors, such as sea surface temperatures, atmospheric pressure systems, and wind shear conditions. The absence of a consistent trend in annual category and maximum winds indicates the inherent variability and unpredictability of hurricane behavior. Understanding these fluctuations provides insights into the diverse range of hurricane impacts and the importance of adaptive disaster management strategies.



The screenshot shows a Google Colab notebook titled "Untitled32.ipynb". The left sidebar displays a file explorer with a folder named "sample_data" containing two files: "cleaned_dataset.csv" and "international_hurricanes.csv". The main code cell contains a list of annual maximum winds in mph, indexed from 0 to 33. The list shows values ranging from 70 to 140 mph. The bottom status bar indicates that the code was completed at 23:22 on 03-03-2024.

```
[33] Annual Maximum Winds (mph):  
Year  
1946.0 80  
1947.0 95  
1948.0 90  
1950.0 90  
1951.0 75  
1952.0 75  
1953.0 75  
1954.0 85  
1955.0 95  
1956.0 75  
1958.0 95  
1960.0 70  
1961.0 75  
1963.0 95  
1964.0 70  
1965.0 105  
1966.0 90  
1967.0 110  
1969.0 95  
1970.0 110  
1983.0 70  
1985.0 95  
1987.0 95  
1988.0 65  
1989.0 120  
1990.0 85  
1992.0 140  
1993.0 85  
1995.0 70  
1996.0 90  
1998.0 70  
1999.0 85  
2000.0 80  
2001.0 80  
2002.0 90  
2003.0 85  
2004.0 95  
2005.0 90  
2007.0 85  
2008.0 85  
2009.0 70  
nan.0 nan
```

A screenshot of a Google Colab notebook titled 'Untitled32.ipynb'. The notebook is open to a code cell containing a list of hurricane data. The data is presented as a list of tuples, where each tuple represents a hurricane. The first element of each tuple is the year, the second is the name, and the third is the maximum wind speed in knots. The data is sorted by year. The notebook interface includes a file explorer on the left showing a directory structure with 'sample_data', 'cleaned_dataset.csv', and 'international_hurricanes.csv'. The bottom status bar shows '0s completed at 23:22' and '2323 03-03-2024'.

Explain your methods and why you chose them

The methods used for this analysis included web scraping with BeautifulSoup to extract data from the NOAA website, data cleaning and manipulation with pandas, and exploratory analysis with pandas functions such as groupby() and aggregate functions. BeautifulSoup was chosen for its ability to parse HTML content and extract tabular data. Pandas was used for its robust data manipulation capabilities, making it easier to clean, filter, and analyze the hurricane dataset. Grouping and aggregating functions helped in summarizing the data, providing insights into trends and patterns. The combination of these tools enabled efficient extraction, cleaning, and analysis of the hurricane data, leading to meaningful interpretations and insights into landfall months, locations, and trends in hurricane activity.