

Configuration Manual

MSc Research Project Cloud Computing

Franklin Ebuka Onyia Student ID: X21221600

School of Computing National College of Ireland

Supervisor:

Rejwanul Haque

National College of Ireland



MSc Project Submission Sheet

School of Computing					
Student Name:					
Student ID:	x21221600				
Programm e:	Msc Cloud Computing Year:2023/2024				
Module:	Research Project				
Lecturer: Submissio	Rejwanul Haque				
n Due Date:					
	Data Exposure Analysis of Misconfigured S3 Buckets: A Quantitative				
Project Title:	Approach				
Word Count:					

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

...Franklin Ebuka Onyia.....

Signature:

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple	
copies)	
Attach a Moodle submission receipt of the online project	
submission, to each project (including multiple copies).	
You must ensure that you retain a HARD COPY of the project,	
both for your own reference and in case a project is lost or mislaid. It is	
not sufficient to keep a copy on computer.	

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

0	office Use	e Only				
••••••			 		 	

Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Franklin Ebuka Onyia Student ID: 21221600

1. Introduction

- **Project Overview**: This manual aids in running code to analyze misconfigured S3 buckets using Python in Google Colab.
- Audience: For anyone that wants to run the analysis.

2. System Requirements

- Hardware: computer with stable internet availability.
- Software: use google chrome Web browser for accessing Google Colab.

3. Environment Setup

- Use Google Colab in your browser.
- Upload the Python notebook file: Click on File > Upload notebook.

4. Configuration of Data

- Follow instructions in the notebook to set up directories and load data.
- Run the blocks of code and setup the data.

5. Execution of Code

- Execute the code blocks in sequence as shown in the notebook.
- Ensure each block runs successfully before moving to the next.

6. Analysis and Visualization of Data

- First Run the data cleaning, transformation, and analysis.
- Visualization tools like Matplotlib and Seaborn are used within the notebook for data representation.

7. Model Evaluation and Regression Analysis

- Follow the code for correlation and regression analysis.
- Evaluate the outputs with respect to the instructions in the notebook.

8. Screenshots and Visual Guide

How to Run Code Blocks Below: click the run button for each step so the code can be ran successfully.

- Extracting and Saving AWS S3 Bucket Information Using Grayhat Warfare API
 API
- [] 1 import requests 2 import json 4 # API key and endpoint setup
 5 api_key = "454286eef0b340f0b6b194152698d786"
 6 url = "https://buckets.grayhatwarfare.com/api/v2/files" 7 headers = {"Authorization": f"Bearer {api_key}"} 8 params = { 9 "type": "aws", # Filter for Amazon S3 buckets 10 "start": 0, # Start from the beginning of the result set 11 "limit": 1000 # Limit of 1000 results 12 } 13 14 # Making the API request 15 response = requests.get(url, headers=headers, params=params) 16 17 # Processing the response 18 if response.status_code == 200: 19 data = response.json() file_info = [] 20 21 21
 # Extracting file information and filtering for AWS buckets
 23
 for file in data.get("files", []):
 24
 if file.get("type") == "aws": # Filter to include only AWS buckets
 25
 | | file_info.append({
- Converting AWS S3 Bucket Data from JSON to CSV Format



Data Preparation

✓ Data Cleaning

```
[ ] 1 import pandas as pd
     2 from datetime import datetime
      2
     4 # Load the dataset
     5 df = pd.read_csv('/content/bucket_file_info.csv')
      6
     7 # Convert 'Last Modified' from Unix timestamp to readable date format
     8 df['Last Modified'] = pd.to_datetime(df['Last Modified'], unit='s')
     9
     10 # Ensure 'Size' is an integer format if it's not already
    11 df['Size'] = df['Size'].astype(int)
    12
    13 # Save the cleaned data back to a new CSV file
    14 df.to_csv('/content/bucket_file_info_cleaned.csv', index=False)
    15
    16 # Display the first few rows to verify
     17 print(df.head())
     18
```

Data Transformation

```
[] 1 import pandas as pd
2
3 # Load the dataset
4 df = pd.read_csv('/content/bucket_file_info_cleaned.csv')
5
6 # Extract file extensions and save in a new column 'File Extension'
7 df['File Extension'] = df['Filename'].str.extract(r'(\.[^.]+)$')
8
9 # Save the updated dataframe back to a new CSV file
10 df.to_csv('/content/bucket_file_info_transformed.csv', index=False)
11
12 # Display the first few rows to verify
13 print(df.head())
14
```

Character of Exposed Data Analysis

Content Type Analysis

```
[] 1 import pandas as pd
     2 import matplotlib.pyplot as plt
     3
     4 # Step 1: Data Loading
     5 df = pd.read_csv('/content/bucket_file_info_transformed.csv')
     6
     7 # Step 2: Data Inspection
     8 # Checking for inconsistencies in 'File Extension' column
     9 if df['File Extension'].isnull().any():
    10 print("There are missing file extensions.")
    11
    12 # Step 3: Data Categorization
    13 # Grouping data by 'File Extension'
    14 file_types = df.groupby('File Extension').size().reset_index(name='Count')
    15
    16 # Step 4: Frequency Calculation
    17 # Sorting the file types by count
    18 file_types = file_types.sort_values('Count', ascending=False)
     19
    20 # Step 5: Data Visualization
     21 # Bar chart to display the frequency distribution of file types
     22 plt.figure(figsize=(10, 8))
     23 plt.bar(file_types['File Extension'], file_types['Count'], color='skyblue')
     24 plt.xlabel('File Extension')
     25 plt.ylabel('Frequency')
     26 plt.title('Frequency Distribution of File Types')
     27 plt.xticks(rotation=45)
     28 plt.tight_layout() # Adjusts plot to ensure everything fits without overlappin
```

Exposure Timeframe Analysis

Last Modified Date Analysis

```
1 import pandas as pd
O
     2 from datetime import datetime
     3
     4 # Load the data
    5 data = pd.read_csv('/content/bucket_file_info_transformed.csv')
     6
    7 # Convert 'Last Modified' to datetime
    8 data['Last Modified'] = pd.to_datetime(data['Last Modified'])
     9
    10 # Calculate the time elapsed since the last modification
    11 current_time = datetime.now()
    12 data['Time Elapsed'] = current_time - data['Last Modified']
    13
    14 # Calculate mean, median, and mode of 'Time Elapsed'
    15 mean_elapsed = data['Time Elapsed'].mean()
    16 median_elapsed = data['Time Elapsed'].median()
    17 mode_elapsed = data['Time Elapsed'].mode()[0] # Taking the first value if there are multiple modes
    18
    19 # Print out the calculated statistics
    20 print(f"Mean Time Elapsed: {mean_elapsed}")
    21 print(f"Median Time Elapsed: {median_elapsed}")
    22 print(f"Mode Time Elapsed: {mode_elapsed}")
    23
    24 # Save the DataFrame with the new calculations
    25 output_file_path = '/content/bucket_file_info_with_elapsed_time.csv'
    26 data.to_csv(output_file_path, index=False)
    27
```

A table summarizing the statistics

```
O
   1 import pandas as pd
     2 import matplotlib.pyplot as plt
     3 from datetime import datetime
     4
     5 # Assuming 'data' is your DataFrame and it already contains 'Time Elapsed' as a timedelta
     6 # Calculate the statistics again
     7 mean_elapsed = data['Time Elapsed'].mean()
     8 median_elapsed = data['Time Elapsed'].median()
     9 mode_elapsed = data['Time Elapsed'].mode()[0]
    10
    11 # Create a summary DataFrame
    12 summary df = pd.DataFrame({
          'Statistic': ['Mean', 'Median', 'Mode'],
    13
    14
          'Time Elapsed': [mean_elapsed, median_elapsed, mode_elapsed]
    15 })
    16
    17 # Display the summary table
    18 print(summary_df)
    19
```

A histogram and a boxplot to visualize the distribution of days since the last modification

```
[] 1 # Histogram of time elapsed
     2 plt.figure(figsize=(10,6))
     3 data['Time Elapsed'].dt.days.plot(kind='hist', bins=30)
     4 plt.title('Histogram of Time Elapsed Since Last Modified')
     5 plt.xlabel('Days')
     6 plt.ylabel('Frequency')
     7 plt.show()
     8
     9 # Boxplot of time elapsed
    10 plt.figure(figsize=(10,6))
    11 data['Time Elapsed'].dt.days.plot(kind='box')
     12 plt.title('Boxplot of Time Elapsed Since Last Modified')
    13 plt.ylabel('Days')
     14 plt.grid(True)
    15 plt.show()
     16
```

Size Analysis

```
1 import pandas as pd
\triangleright
    2 import matplotlib.pyplot as plt
    3 from scipy.stats import pearsonr
    4 from datetime import datetime
    5
    6 # Load the dataset
    7 df = pd.read_csv('/content/bucket_file_info_transformed.csv')
    8
    9 df['Last Modified'] = pd.to_datetime(df['Last Modified'], format='%Y-%m-%d %H:%M:%S')
   10
   11
   12 # Calculate the number of days since the oldest file for each entry
   13 oldest_date = df['Last Modified'].min()
   14 df['Days Since Oldest'] = (df['Last Modified'] - oldest_date).dt.days
   15
   16 # Summary statistics for file sizes
   17 size_stats = df['Size'].describe()
   18
   19 # Correlational analysis
   20 correlation = pearsonr(df['Size'], df['Days Since Oldest'])[0]
   21
   22 # Create scatter plot for the correlational analysis
   23 plt.figure(figsize=(10, 6))
   24 plt.scatter(df['Days Since Oldest'], df['Size'])
   25 plt.title('File Size vs Days Since Oldest File')
   26 plt.xlabel('Days Since Oldest File')
   27 plt.ylabel('File Size (bytes)')
    28 plt.yscale('log') # Using a log scale due to wide range of file sizes
    29 plt.grid(True)
```

Correlation and Regression Analysis

```
1 import pandas as pd
 2 import numpy as np
 3 import matplotlib.pyplot as plt
 4 import seaborn as sns
 5 from sklearn.linear_model import LinearRegression
 6 from sklearn.model selection import train test split
 7 from sklearn.metrics import mean_squared_error, r2_score
 8
9 # Load the dataset
10 df = pd.read_csv('/content/bucket_file_info_transformed.csv')
11
12 # Convert 'Last Modified' to datetime and 'Size' to numeric if necessary
13 df['Last Modified'] = pd.to_datetime(df['Last Modified'])
14 df['Size'] = pd.to_numeric(df['Size'])
15
16 # Calculate 'Exposure Time' in days
17 df['Exposure Time'] = (pd.to_datetime('today') - df['Last Modified']).dt.days
18
19 # EDA: Basic stats and visualizations
20 print(df.describe())
21 sns.pairplot(df[['Size', 'Exposure Time']])
22 plt.show()
23
24 # Correlation Analysis
25 correlation_matrix = df.corr()
26 print(correlation_matrix)
```

9. Troubleshooting and Common Issues

- In case of any problem, please re-check code syntax, also make sure all libraries are installed properly.
- In case of any issue, please, checkout Google Colab's help resources or Python programming learning sites.

10. Security Considerations

- Be sure you use libraries that are from their official website.
- Also be careful with Security and privacy issues, most especially when you have used sensitive information.