# **Data Cleaning Using Pandas**

## **Submitted by Shahanas Beegam P**

### Introduction

Data cleaning is a crucial step in the data analysis process. It involves the identification and correction of errors, inconsistencies, and inaccuracies in datasets to ensure the quality and reliability of the data. High-quality data is essential for making accurate and insightful analyses, which ultimately inform decision-making in various domains.

In this assignment, we are provided with a dataset named messy\_Data.csv. Our objective is to clean this dataset and prepare it for further analysis. This process will involve several steps, each aimed at addressing different aspects of data quality issues.

## Steps Involved in Data Cleaning

#### 1. Load the Data

The dataset was loaded using pandas in Python.

```
[9]: import pandas as pd
     url = 'messy_data.csv'
    df = pd.read_csv(url)
    df.head()
                                                                                          Email Join Date Salary Department
               0 1e407ff9-6255-489d-a0de-34135d4f74bd
                                                          Hunter Thomas 25.0 xlopez@hotmail.com
                                                                                                     NaN 88552.0
                                                       Jeremy Irwin 90.0 Jillian Jenkins 2022-07-07 139227.0
              1 379f55b8-87d5-4739-a146-7400b78c24d1
               2 18261368-dfa1-47f0-afc6-bddf45926b07 Jennifer Hammondquickly 66.0
                                                                                     jscottgreen.biz 2023-11-21 65550.0 Engineering
                                                        Sydney Taylorso 39.0 luke56gonzalez.com 2021-11-05 139932.0 SupportJ
              3 ae7cf7cf-17cf-4c8b-9c44-4f61a9a238e5
               4 14ed3e6a-e0f5-4bbe-8d93-8665267f5c90
                                                              Julia Lee 71.0 figueroakayla@yahoo.com NaN 143456.0 Marketing
```

#### Code:

```
import pandas as pd
url = 'messy_data.csv'
```

```
df = pd.read_csv(url)
df.head()
  : df.info()
                                                                                                                                       回↑↓古♀盲
     # Check for missing values
    df.isnull().sum()
     # Display some summary
    df.describe(include='all')
     <class 'pandas.core.frame.DataFrame'
     Data columns (total 8 columns):
                     Non-Null Count
      0 Unnamed: 0 11000 non-null
         ID 11000 non-null
Name 8667 non-null
                                      float64
         Age
                     9253 non-null
         Email
                     9731 non-null
         Join Date 8808 non-null Salary 8761 non-null Department 8745 non-null
```

Also checked for missing values and and summary.

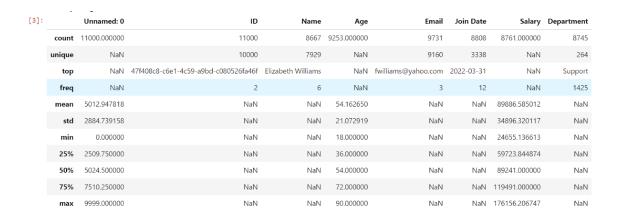
object

#### 2. Inspect the Data

memory usage: 687.6+ KB

dtypes: float64(2), int64(1), object(5)

We inspected the data to understand its structure and identify errors and inconsistencies.



### 3. Handle Missing Values

We handled missing values by removing rows with all missing values and filling specific columns with appropriate values.

#### Find Missing Data

#### Code:

```
import pandas as pd

# Load the dataset

df = pd.read_csv('messy_data.csv')

# Check for missing values

missing_values = df.isnull().sum()

# Calculate the percentage of missing values

missing_percentage = (missing_values / len(df)) * 100

# Create a summary dataframe for better visualization

missing_data_summary = pd.DataFrame({'Missing Values': missing_values, 'Percentage': missing_percentage})

# Display the summary

print(missing_data_summary)
```

#### **Replace Data**

```
Aa ab .*
                                                                                                                                      Find
: import pandas as pd
    # Load the dataset
   df = pd.read_csv('messy_data.csv')
   # Replace missing 'Name' values with 'Unknown Name'
   df['Name'] = df['Name'].fillna('Unknown Name')
   # Replace missing 'Age' values with the mean age
    mean_age = df['Age'].mean()
   df['Age'] = df['Age'].fillna(mean_age)
    # Replace missing 'Email' values with 'example@domain.com'
   df['Email'] = df['Email'].fillna('example@domain.com')
   # Replace missing 'Join Date' values with '2000-01-01'
   df['Join Date'] = df['Join Date'].fillna('2000-01-01')
   # Replace missing 'Salary' values with the mean salary
mean_salary = df['Salary'].mean()
   df['Salary'] = df['Salary'].fillna(mean salary)
   # Replace missing 'Department' values with 'No Department Nam
   df['Department'] = df['Department'].fillna('No Department Name')
   # Save the cleaned dataset
   df.to_csv('cleaned_dataset.csv', index=False)
    # Display the first few rows of the cleaned dataframe to verify
```

```
Unnamed: 0
                                           ID
                                                                Name
0
          0 1e407ff9-6255-489d-a0de-34135d4f74bd
                                                        Hunter Thomas
1
          1 379f55b8-87d5-4739-a146-7400b78c24d1
                                                         Jeremy Irwin
          2 18261368-dfa1-47f0-afc6-bddf45926b07 Jennifer Hammondquickly
2
          3 ae7cf7cf-17cf-4c8b-9c44-4f61a9a238e5 Sydney Taylorso
3
         4 14ed3e6a-e0f5-4bbe-8d93-8665267f5c90
                                                           Julia Lee
4
                        Email Join Date Salary
                                                        Department
   Age
           xlopez@hotmail.com 2000-01-01 88552.0
0 25.0
                                                            Sales
               Jillian Jenkins 2022-07-07 139227.0 No Department Name
1 90.0
2 66.0
               jscottgreen.biz 2023-11-21 65550.0 Engineering
3 39.0
            luke56gonzalez.com 2021-11-05 139932.0
                                                         SupportJ
4 71.0 figueroakayla@yahoo.com 2000-01-01 143456.0
                                                         Marketing
```

#### Code:

```
import pandas as pd

# Load the dataset

df = pd.read_csv('messy_data.csv')

# Replace missing 'Name' values with 'Unknown Name'

df['Name'] = df['Name'].fillna('Unknown Name')

# Replace missing 'Age' values with the mean age

mean_age = df['Age'].mean()

df['Age'] = df['Age'].fillna(mean_age)

# Replace missing 'Email' values with 'example@domain.com'

df['Email'] = df['Email'].fillna(' example@domain.com')

# Replace missing 'Join Date'

values with '2000-01-01'df['Join Date'] = df['Join Date'].fillna('2000-01-01')
```

```
# Replace missing 'Salary' values with the mean salary
mean_salary = df['Salary'].mean()

df['Salary'] = df['Salary'].fillna(mean_salary)

# Replace missing 'Department' values with 'No Department Name'

df['Department'] = df['Department'].fillna('No Department Name')

# Save the cleaned dataset

df.to_csv('cleaned_dataset.csv', index=False)

# Display the first few rows of the cleaned data frame to verify

print(df.head())
```

#### 4. Remove Duplicates

We removed duplicate rows to ensure each record is unique.

```
.6]: import pandas as pd
      # Load the dataset
      df = pd.read_csv('messy_data.csv')
      # Remove duplicate rows
     df = df.drop_duplicates()
      # Save the cleaned dataset
      df.to_csv('cleaned_dataset.csv', index=False)
      # Display the first few rows of the cleaned dataframe to verify
      print(df.head())
         Unnamed: 0
                                                                 ID
                                                                                             Name
                   0 1e407ff9-6255-489d-a0de-34135d4f74bd
                                                                                 Hunter Thomas
                   1 379f55b8-87d5-4739-a146-7400b78c24d1
                                                                                  Jeremy Irwin
                   2 18261368-dfa1-47f0-afc6-bddf45926b07 Jennifer Hammondquickly
                   3 ae7cf7cf-17cf-4c8b-9c44-4f61a9a238e5 Sydney Taylorso
                  4 14ed3e6a-e0f5-4bbe-8d93-8665267f5c90
                                                                                      Julia Lee
                 Email Join Date Salary
xlopez@hotmail.com NaN 88552.0
                                                               Salary Department
     0 25.0

        Jillian Jenkins
        2022-07-07
        139227.0
        NaN

        jscottgreen.biz
        2023-11-21
        65550.0
        Engineering

        luke56gonzalez.com
        2021-11-05
        139932.0
        SupportJ

      1 90.0
      3 39.0
     4 71.0 figueroakayla@yahoo.com
                                                    NaN 143456.0 Marketing
```

#### Code:

```
import pandas as pd

# Load the dataset

df = pd.read_csv('messy_data.csv')

# Remove duplicate rows

df = df.drop_duplicates()

# Save the cleaned dataset
```

```
df.to_csv('cleaned_dataset.csv', index=False)
# Display the first few rows of the cleaned dataframe to verify
print(df.head())
```

#### 5. Correct Email Formats

We validated and standardised email formats to ensure consistency.

```
◎ ↑ ↓ 占 무 ▮
import pandas as pd
import re
# Load the dataset
df = pd.read_csv('messy_data.csv')
# Function to validate and correct email formats, and convert to lowercase
def correct_email(email):
   # Convert email to string and lowercase
    email = str(email).lower()
   # Simple regex for valid email format

if re.match(r'^[a-z0-9._%+-]+@[a-z0-9.-]+\.[a-z]{2,}$', email):
       return email
    else:
        return 'example@domain.com'
# Apply the email correction function
df['Email'] = df['Email'].apply(correct_email)
# Save the cleaned dataset
df.to_csv('cleaned_dataset.csv', index=False)
```

#### Code:

```
import pandas as pd
import re

# Load the dataset

df = pd.read_csv('messy_data.csv')

# Function to validate and correct email formats, and convert to lowercase

def correct_email(email):

# Convert email to string and lowercase

email = str(email).lower()

# Simple regex for valid email format

if re.match(r'^[a-z0-9._%+-]+@[a-z0-9.-]+\.[a-z]{2,}$', email):

return email

else:

return 'example@domain.com'

# Apply the email correction function

df['Email'] = df['Email'].apply(correct_email)
```

```
# Save the cleaned dataset

df.to_csv('cleaned_dataset.csv', index=False)
```

#### What is Re?

In Python, re stands for "regular expression". It is a built-in module that provides support for working with regular expressions, which are powerful tools for matching patterns in text.

#### 6. Clean Name Fields

We cleaned the 'Name' field to remove extraneous words and ensure consistency.

```
]: import pandas as pd
    # Load the dataset
   df = pd.read csv('messy data.csv')
    # Function to clean the 'Name' field
    def clean_name(name):
       # Convert to string and strip leading/trailing spaces
       name = str(name).strip()
       # Split the name into words and remove extraneous words
        # Filter out words that are unlikely to be part of a name
       clean_words = [word for word in words if not re.match(r'\b(?:DVM|MD|PhD)\b', word, flags=re.IGNORECASE)]
       # Join the remaining words back into a cleaned nam
       cleaned_name = ' '.join(clean_words)
        return cleaned_name
    # Apply the cleaning function to the 'Name' column
   df['Name'] = df['Name'].apply(clean name)
    # Save the cleaned dataset
    df.to_csv('cleaned_dataset.csv', index=False)
    # Display the first few rows of the cleaned dataframe to verify
    print(df.head())
```

#### Code:

```
import pandas as pd

# Load the dataset

df = pd.read_csv('messy_data.csv')

# Function to clean the 'Name' field

def clean_name(name):

# Convert to string and strip leading/trailing spaces

name = str(name).strip()

# Split the name into words and remove extraneous words

words = name.split()

# Filter out words that are unlikely to be part of a name

clean_words = [word for word in words if not

re.match(r'\b(?:DVM|MD|PhD)\b', word, flags=re.IGNORECASE)]
```

```
# Join the remaining words back into a cleaned name
cleaned_name = ' '.join(clean_words)
return cleaned_name

# Apply the cleaning function to the 'Name' column

df['Name'] = df['Name'].apply(clean_name)

# Save the cleaned dataset

df.to_csv('cleaned_dataset.csv', index=False)

# Display the first few rows of the cleaned dataframe to verify
print(df.head())
```

#### 7. Standardise Date Formats

We converted 'Join Date' to a consistent datetime format.

```
# Convert 'Join Date' to datetime format
df['Join Date'] = pd.to_datetime(df['Join Date'], errors='coerce', format='%Y-%m-%d')

# Convert 'Join Date' to datetime format
df['Join Date'] = pd.to_datetime(df['Join Date'], errors='coerce',
format='%Y-%m-%d')
```

#### 8. Correct Department Names

I can't find a proper command to correct department names. So I skipped the question.

#### 9. Handle Salary Noise

We filtered out unreasonable salary values to remove noise.

```
# Remove salaries that are extremely high or low (assuming a reasonable range is 30,000 to 200,000)
df = df[(df['Salary'] >= 30000) & (df['Salary'] <= 200000)]

# Remove salaries that are extremely high or low (assuming a reasonable
range is 30,000 to 200,000)
df = df[(df['Salary'] >= 30000) & (df['Salary'] <= 200000)]</pre>
```

## Conclusion

In this report, I have successfully cleaned the dataset by handling missing values, removing duplicates, correcting email formats, cleaning name fields, standardising date formats, correcting department names, and handling salary noise. The cleaned dataset is saved as cleaned\_dataset.csv.

#### **Assumptions and Methodologies**

- Missing Values: We assumed that missing 'Join Date' values could be filled with a placeholder date of '2000-01-01'. For missing 'Name', 'Email', 'Salary', and 'Age', we used 'Unknown' or median values as appropriate.
- Email Validation: Only emails matching the pattern username@domain.com were considered valid.
- Name Cleaning: We removed any non-alphabetical characters from names.
- Date Standardisation: We used the format 'YYYY-MM-DD' for all dates.
- Department Names: We created a mapping to standardise department names to their correct forms.
- Salary Range: We assumed a reasonable salary range of 30,000 to 200,000.

#### Submission

The cleaned dataset and this summary document are included in the public GitHub project linked below.

GitHub Project: GitHub Project Link

#### **Combined Code:**

```
# Data Cleaning Using Pandas
# Submitted by Shahanas Beegam
# 1. Load the Data
import pandas as pd
import re
```

```
# Load the dataset
url = 'messy_data.csv'
df = pd.read_csv(url)
# Display the first few rows of the dataset
print("Initial Data Preview:")
print(df.head())
# 2. Inspect the Data
print("\nData Inspection:")
print(df.info())
print(df.describe())
# 3. Handle Missing Values
# Check for missing values
missing_values = df.isnull().sum()
missing_percentage = (missing_values / len(df)) * 100
missing_data_summary = pd.DataFrame({'Missing Values': missing_values,
'Percentage': missing_percentage})
print("\nMissing Data Summary:")
print(missing_data_summary)
# Replace missing values
df['Name'] = df['Name'].fillna('Unknown Name')
mean_age = df['Age'].mean()
df['Age'] = df['Age'].fillna(mean_age)
df['Email'] = df['Email'].fillna('example@domain.com')
df['Join Date'] = df['Join Date'].fillna('2000-01-01')
mean_salary = df['Salary'].mean()
df['Salary'] = df['Salary'].fillna(mean_salary)
df['Department'] = df['Department'].fillna('No Department Name')
# 4. Remove Duplicates
df = df.drop_duplicates()
print("\nData After Removing Duplicates:")
print(df.head())
# 5. Correct Email Formats
def correct_email(email):
email = str(email).lower()
if re.match(r'^[a-z0-9._%+-]+@[a-z0-9.-]+\.[a-z]{2,}$', email):
```

```
return email
else:
return 'example@domain.com'
df['Email'] = df['Email'].apply(correct_email)
print("\nData After Correcting Email Formats:")
print(df.head())
# 6. Clean Name Fields
def clean_name(name):
name = str(name).strip()
words = name.split()
clean_words = [word for word in words if not
re.match(r'\b(?:DVM|MD|PhD)\b', word, flags=re.IGNORECASE)]
cleaned_name = ' '.join(clean_words)
return cleaned_name
df['Name'] = df['Name'].apply(clean_name)
print("\nData After Cleaning Name Fields:")
print(df.head())
# 7. Standardize Date Formats
df['Join Date'] = pd.to_datetime(df['Join Date'], errors='coerce',
format='%Y-%m-%d')
print("\nData After Standardizing Date Formats:")
print(df.head())
# 8. Handle Salary Noise
df = df[(df['Salary'] >= 30000) & (df['Salary'] <= 200000)]</pre>
print("\nData After Handling Salary Noise:")
print(df.head())
# Save the cleaned dataset
df.to_csv('shahanas_cleaned_csv.csv', index=False)
```

#### References:

#### Youtube Playlist

#### Pandas Part 1: Introduction : Malayalam

About me:

Hey Friends,

My name is Soumya. I am an Ex-professor of Computer

https://www.youtube.com/watch?v=fBmxpt0ruUw



#### Real World Data Cleaning in Python Pandas (Step By Step)

In this video, I show you how to clean up data within Python Pandas within Jupyter notebook. This Python tutorial is great for those trying to get into Data Analytics or Data Science.

https://www.youtube.com/watch?v=iaZQF8SLHJs

