

Assignment-1

Data Cleaning (Pandas)

By [Arjit Dhall](#) [Linkedin](#)

This dataset helps you to rehearse the data-cleaning process using the pure Python pandas library.

Columns:

- Age
- Salary
- Rating
- Location
- Established
- Easy Apply

Here are some questions that you can follow to perform data cleaning and manipulation using the given dataset.

The actual raw.csv data is given below

```
In [ ] : import pandas as pd
url='https://raw.githubusercontent.com/ArjitDhall/PrepInsta-DA-Week-3/main/Assignment1.csv'
df=pd.read_csv(url,encoding='unicode_escape')
df

Out [ ] :
   Index  Age  Salary  Rating  Location  Established  Easy Apply
0      0  4.0  44k-59k  5.4  India,in  1999  TRUE
1      1  6.0  55k-69k  5.3  New York,Ny  2002  TRUE
2      2  3.0  77k-99k  1.0  Australia  2000  TRUE
3      3  4.0  44k-59k  4.4  India,in  1999  TRUE
4      4  25.0  44k-59k  6.4  Australia  2002  -1
5      5  44.0  77k-99k  1.4  India,in  1999  TRUE
6      6  21.0  44k-59k  0.0  New York,Ny  -1  -1
7      7  44.0  44k-59k  -1.0  Australia  -1  -1
8      8  35.0  44k-59k  5.4  New York,Ny  -1  -1
9      9  22.0  44k-59k  7.7  India,in  -1  TRUE
10     10  55.0  10k-49k  5.4  India,in  2008  TRUE
11     11  44.0  10k-49k  0.0  India,in  2009  -1
12     12  NaN  44k-59k  6.7  India,in  2009  -1
13     13  25.0  44k-59k  -1.0  Australia  2019  TRUE
14     14  60.0  44k-59k  4.0  Australia  2020  TRUE
15     15  44.0  88k-100k  3.0  Australia  1999  -1
16     16  19.0  19k-40k  4.5  India,in  1984  -1
17     17  NaN  44k-59k  5.3  New York,Ny  1943  TRUE
18     18  35.0  44k-59k  6.7  New York,Ny  1954  TRUE
19     19  32.0  44k-59k  3.3  New York,Ny  1955  TRUE
20     20  NaN  44k-59k  5.7  New York,Ny  1944  TRUE
21     21  35.0  44k-59k  5.0  New York,Ny  1946  -1
22     22  18.0  55k-69k  7.8  New York,Ny  1988  TRUE
23     23  20.0  44k-59k  2.0  New York,Ny  1999  TRUE
24     24  20.0  44k-59k  -1.0  New York,Ny  1987  -1
25     25  55.0  44k-59k  0.0  Australia  1980  TRUE
26     26  NaN  55k-69k  NaN  India,in  1934  TRUE
27     27  52.0  44k-59k  5.4  India,in  1935  -1
28     28  NaN  39k-49k  5.4  Australia  1932  -1
```

Problem Solutions:

1. Expand to reveal solutions to the problems. [Orderwise]

1. Missing Values:

Question:

Are there any missing values in the dataset, and if so, how should they be handled for each indicator?

Answer:

- Yes: There are missing values in dataset. There are missing values in multiple columns (e.g., Age, Rating, Established, Easy Apply).
- Handling: These missing values need to be addressed, either through imputation, removal of rows, or depending on the context.
- However: column named "Easy Apply" has -1 which means False
- Whereas for columns "Age", "Rating", "Established", we need to replace with mean of Age, mean of Rating and "Unknown" values respectively.

The solution to clean missing values:

```
df["Easy Apply"] = df["Easy Apply"].str.replace('-1', 'FALSE') # Replace -1 with False in "Easy Apply Column"

age_mean = df["Age"].mean() # For Replacing NaN values with mean of Age
age_mean = int(age_mean)
df["Age"] = df["Age"].replace(np.nan, age_mean)

rating_mean = df["Rating"].mean() # For Replacing NaN value with mean of Rating
rating_mean = int(rating_mean)
df["Rating"] = df["Rating"].replace(np.nan, rating_mean)

df["Established"] = df["Established"].replace('-1', 'Unknown') # Replace -1 with Unknown
```

2. Data Types:

Question:

What are the data types of each indicator, and do they align with their expected types (e.g., numerical, categorical)?

Answer:

- Age: Mostly numerical (except for some missing values).
- Salary: String format (needs conversion to numerical values).
- Rating: Floating numerical with some missing values represented as -1.
- Location: String format.
- Established: Mostly numerical with some missing values.
- Easy Apply: Boolean (except for missing values).

The solution to find out the data type of each series:

```
df[["Indicator/Column_Name"]].dtype # Returns the fdata type of each indicator
```

3. Outliers:

Identify potential outliers in numerical indicators (e.g., Age, Salary, Rating). Should outliers be removed or adjusted?

Answer:

- Age: Any age that seems unusually high or low compared to the typical range in the dataset could be considered an outlier. For example, 13 is seen to be low outlier and 86 seems to be high outlier.
- Salary: Outliers might include exceptionally high or low salaries compared to the majority of salaries in the dataset. For example, 100000 is seen to be low outlier and 1000000 seems to be high outlier Salary.
- Rating: Values significantly deviating from the general range of ratings could be identified as outliers. For example, 8 is seen to be low outlier and 7.7 seems to be high outlier Rating.

The solution to find out the outliers, we use the following algorithm:

```
q1 = data.quantile(0.25)
q3 = data.quantile(0.75)
iqr = q3 - q1
lower_bound = q1 - 1.5 * iqr
upper_bound = q3 + 1.5 * iqr
outliers = data[(data < lower_bound) | (data > upper_bound)]
```

4. Salary Formatting:

Question:

Examine the format of the Salary column. Does it require any formatting or standardization for consistent analysis?

Answer:

- Salary values are in the format "\$XXX-\$XXX" and require transformation into numerical values.

- Split the "Salary" column at "-", forming two columns namely "Starting Salary" and "Ending Salary"
- Left Strip the '\$' symbol for both of the columns.
- Replace 'k' alphabet to '000'

Following the below snippet, we can clean the "Salary" column

```
df[["Starting Salary", "Ending Salary"]] = df["Salary"].str.split('-', 1, expand=True) # Splits the following Salary column
df.drop(columns=["Salary"]) # Delete the Salary column

df["Starting Salary"] = df["Starting Salary"].str.strip("$") # Left Stripping '$'
df["Ending Salary"] = df["Ending Salary"].str.strip("$") # Mapped the string to respective
df["Starting Salary"] = df["Starting Salary"].str.replace('k', '000') # Replace 'k' with '000'
df["Ending Salary"] = df["Ending Salary"].str.replace('k', '000')
```

5. Location Standardization:

Question:

Check the consistency of location entries. Do they need standardization, and how can this be achieved?

Answer:

- Locations have various formats (e.g., "India,in", "New York,Ny", "Australia,Aus"). Standardization could involve splitting these into separate columns for "Country" and "Initials". Let us take a tour of all the countries.
- India,in : Here we have to split India and In at ','.
- New York,Ny : Here also we need to split New York and Ny at ','.
- Australia,Aus : Its a special case where we need to split Australia and Aus at ',' but it shouldn't collide with New York, so we need to place them prioritywise.

The following steps show how to split the Location into two separate columns.

```
# Splitting based on different separators for different locations
def split_location(row):
    location = row["Location"]
    if ',' in location:
        parts = location.split(',')
        country = parts[0].strip()
        initials = parts[1].strip()
    elif '-' in location:
        parts = location.split('-')
        country = parts[0].strip()
        initials = parts[1].strip()
    else:
        country = location.strip()
        initials = None # Or set initials as needed for cases with no separator
    return pd.Series({'Country': country, 'Initials': initials})

# Apply the function to split 'Location' into 'Country' and 'Initials'
df[["Country", "Initials"]] = df.apply(split_location, axis=1)
df.drop(columns=["Location"])
```

6. Established Column:

Question:

Explore the Established column. Are there any inconsistencies or anomalies that need to be addressed?

Answer:

Yes, Inconsistencies exist in the Established column, with missing values represented as -1. It is converted to string first, then replace -1 with 'Unknown' string

It can be solved if emphasize the following snippet:

```
df["Established"] = df["Established"].astype(str) # Converting them to string data type
df["Established"] = df["Established"].replace('-1', 'Unknown') # Replace -1 with Unknown
```

7. Easy Apply Indicator:

Question:

Analyze the Easy Apply column. Does it contain boolean values or need transformation for better analysis?

Answer:

Easy Apply: column initially consists of string data type. Apart from that, it also consists of missing values which is "-1". However, we consider -1 as FALSE and replace all -1 with FALSE.

Still, it is carries, string data type, which is converted to bool using map function.

Following algorithm is used for that case:

```
df["Easy Apply"] = df["Easy Apply"].str.replace('-1', 'FALSE') # Replace -1 with False
mapping = {'TRUE': True, 'FALSE': False} # Mapped the string to respective
# bool value
df["Easy Apply"] = df["Easy Apply"].map(mapping).fillna(False).astype(bool) # Formatting the column with bool
lean data type
```

8. Rating Range:

Question:

Investigate the range of values in the Rating column. Does it fall within expected rating scales, and how should outliers be treated?

Answer:

Ratings range from 0 to 7.8, with missing values represented as -1. It is justified that, it lies on a rating scale of 10 points. Moreover, missing values may need addressing.

For that kind of scenarios, we first replace the missing values with NaN then find the mean of "Rating" column, then replace all NaN values with its mean:

```
rating_mean = df["Rating"].mean() # For Replacing NaN value with mean of Rating
rating_mean = int(rating_mean)
df["Rating"] = df["Rating"].replace(np.nan, rating_mean)
```

9. Age Distribution:

Question:

Check the distribution of values in the Age column. Are there any unusual entries, and how might they impact analysis?

Answer:

Age values are present, but some are missing. Apart from that, we can see certain outliers in the Age column. Any age that seems unusually high or low compared to the typical range in the dataset could be considered an outlier. For example, 13 is seen to be low outlier and 86 seems to be high outlier Age. In this way we delete rows where Age is NaN:

```
df = df.dropna() # Delete rows consisting NaN in data frame
```

10. Handling Special Characters:

Question:

Examine all text-based columns (e.g., Location). Are there special characters or inconsistencies that need cleaning?

Answer:

- Locations entries contain commas and spaces that require splitting or cleaning. Locations have various formats (e.g., "India,in", "New York,Ny", "Australia,Aus"). Standardization could involve splitting these into separate columns for "Country" and "Initials". Let us take a tour of all the countries.
- India,in : Here we have to split India and In at ','.
- New York,Ny : Here also we need to split New York and Ny at ','.
- Australia,Aus : Its a special case where we need to split Australia and Aus at ',' but it shouldn't collide with New York, so we need to place them prioritywise.

The following steps show how to split the Location into two separate columns.

```
# Splitting based on different separators for different locations
def split_location(row):
    location = row["Location"]
    if ',' in location:
        parts = location.split(',')
        country = parts[0].strip()
        initials = parts[1].strip()
    elif '-' in location:
        parts = location.split('-')
        country = parts[0].strip()
        initials = parts[1].strip()
    else:
        country = location.strip()
        initials = None # Or set initials as needed for cases with no separator
    return pd.Series({'Country': country, 'Initials': initials})

# Apply the function to split 'Location' into 'Country' and 'Initials'
df[["Country", "Initials"]] = df.apply(split_location, axis=1)
df.drop(columns=["Location"])
```

Salary values are in the format "\$XXX-\$XXX" and require transformation into numerical values.

- Split the "Salary" column at "-", forming two columns namely "Starting Salary" and "Ending Salary"
- Left Strip the '\$' symbol for both of the columns.
- Replace 'k' alphabet to '000'

Following the below snippet, we can clean the "Salary" column

```
df[["Starting Salary", "Ending Salary"]] = df["Salary"].str.split('-', 1, expand=True) # Splits the following Salary column
df.drop(columns=["Salary"]) # Delete the Salary column

df["Starting Salary"] = df["Starting Salary"].str.strip("$") # Left Stripping '$'
df["Ending Salary"] = df["Ending Salary"].str.strip("$") # Mapped the string to respective
df["Starting Salary"] = df["Starting Salary"].str.replace('k', '000') # Replace 'k' with '000'
df["Ending Salary"] = df["Ending Salary"].str.replace('k', '000')
```

11. Data Integrity:

Question:

Ensure data integrity by cross-referencing entries. For instance, does the Established column align with the Age column?

Answer:

Yes, Inconsistencies exist between the Established and Age columns, with missing values causing misalignment.

- The Established column denotes the year of establishment for certain records, represented mostly in numerical format but contains missing values indicated by -1.
- The Age column refers to the age of individuals or entities, but it also contains missing values.
- We can compare the Established columns values (years) with the Age column to identify inconsistencies or misalignments. If both columns refer to the same entity, the established year should correlate with the age or close to it.

12. Easy Apply Transformation:

Question:

If the Easy Apply column contains non-boolean values, how can it be transformed into a usable format?

Answer:

Some missing values exist in the Easy Apply column that may require handling. Easy Apply column initially consists of string data type. Apart from that, it also consists of missing values which is "-1". However, we consider -1 as FALSE and replace all -1 with FALSE.

Still, it is carries, string data type, which is converted to bool using map function.

Following algorithm is used for that case:

```
df["Easy Apply"] = df["Easy Apply"].str.replace('-1', 'FALSE') # Replace -1 with False
mapping = {'TRUE': True, 'FALSE': False} # Mapped the string to respective
# bool value
df["Easy Apply"] = df["Easy Apply"].map(mapping).fillna(False).astype(bool) # Formatting the column with bool
lean data type
```

13. Location Accuracy:

Question:

Assess the accuracy of location entries. Are there misspelled or ambiguous locations that require correction?

Answer:

- Entries like "India,in" and "New York,Ny" could be used in splitting into separate columns for country and initials at ','.
- Apart from that "Australia,Aus" might also be used in splitting into the same two columns at ','.
- Also, we need to look out for that if we split at ',' first then, New York, would be affected. To avoid it we can use them in df["Country", "Initials"] but it shouldn't collide with New York, so we need to place them prioritywise.

```
# Splitting based on different separators for different locations
def split_location(row):
    location = row["Location"]
    if ',' in location:
        parts = location.split(',')
        country = parts[0].strip()
        initials = parts[1].strip()
    elif '-' in location:
        parts = location.split('-')
        country = parts[0].strip()
        initials = parts[1].strip()
    else:
        country = location.strip()
        initials = None # Or set initials as needed for cases with no separator
    return pd.Series({'Country': country, 'Initials': initials})

# Apply the function to split 'Location' into 'Country' and 'Initials'
df[["Country", "Initials"]] = df.apply(split_location, axis=1)
df.drop(columns=["Location"])
df.reset_index(drop=True, inplace=True)
df.drop(columns=["Index"], inplace=True)
df
```

Out [] :

```
   Age  Rating  Established  Easy Apply  Starting Salary  Ending Salary  Country  Initials
0  44  5.4      1999        True      44000      99000      India      in
1  66  3.5      2002        True      55000      66000      New York   Ny
2  39  4.0      Unknown      False      77000      89000      New York   Ny
3  64  4.4      1988        False      44000      99000      India      in
4  25  6.4      2002        False      44000      99000      Australia  Aus
5  44  1.4      1999        True      77000      89000      India      in
6  21  0.0      Unknown      False      44000      99000      New York   Ny
7  44  4.0      Unknown      False      44000      99000      Australia  Aus
8  35  5.4      Unknown      True      44000      99000      New York   Ny
9  22  7.7      Unknown      True      44000      99000      India      in
10  55  5.4      2008        True      10000      49000      India      in
11  44  6.7      2009        False      10000      49000      India      in
12  39  0.0      1999        False      44000      99000      India      in
13  25  4.0      2019        True      44000      99000      Australia  Aus
14  66  4.0      2020        True      44000      99000      Australia  Aus
15  44  3.0      1999        False      89000      101000     Australia  Aus
16  19  4.5      1984        False      19000      40000      India      in
17  39  5.3      1943        True      44000      99000      New York   Ny
18  35  6.7      1954        True      44000      99000      New York   Ny
19  32  3.3      1955        True      44000      99000      New York   Ny
20  39  5.7      1944        True      44000      99000      New York   Ny
21  35  5.0      1946        False      44000      99000      New York   Ny
22  39  7.8      1988        True      55000      66000      New York   Ny
23  39  2.4      1980        True      44000      99000      New York   Ny
24  13  4.0      1987        False      44000      99000      New York   Ny
25  55  0.0      1980        True      44000      99000      Australia  Aus
26  39  4.0      1934        True      55000      66000      India      in
27  52  5.4      1935        False      44000      99000      India      in
28  39  3.4      1932        False      39000      88000      Australia  Aus
```

With a column on top, we will lowercase the Column Headings (It is determined to do at previous step)

```
In [ ] : df.columns = df.columns.str.lower().str.replace(' ', '_')
df

Out [ ] :
   age  rating  established  easy_apply  starting_salary  ending_salary  country  initials
0  44  5.4      1999        True      44000      99000      india      in
1  66  3.5      2002        True      55000      66000      new york   ny
2  39  4.0      unknown      False      77000      89000      new york   ny
3  64  4.4      1988        False      44000      99000      india      in
4  25  6.4      2002        False      44000      99000      australia  aus
5  44  1.4      1999        True      77000      89000      india      in
6  21  0.0      unknown      False      44000      99000      new york   ny
7  44  4.0      unknown      False      44000      99000      australia  aus
8  35  5.4      unknown      True      44000      99000      new york   ny
9  22  7.7      unknown      True      44000      99000      india      in
10  55  5.4      2008        True      10000      49000      india      in
11  44  6.7      2009        False      10000      49000      india      in
12  39  0.0      1999        False      44000      99000      india      in
13  25  4.0      2019        True      44000      99000      australia  aus
14  66  4.0      2020        True      44000      99000      australia  aus
15  44  3.0      1999        False      89000      101000     australia  aus
16  19  4.5      1984        False      19000      40000      india      in
17  39  5.3      1943        True      44000      99000      new york   ny
18  35  6.7      1954        True      44000      99000      new york   ny
19  32  3.3      1955        True      44000      99000      new york   ny
20  39  5.7      1944        True      44000      99000      new york   ny
21  35  5.0      1946        False      44000      99000      new york   ny
22  39  7.8      1988        True      55000      66000      new york   ny
23  39  2.4      1980        True      44000      99000      new york   ny
24  13  4.0      1987        False      44000      99000      new york   ny
25  55  0.0      1980        True      44000      99000      australia  aus
26  39  4.0      1934        True      55000      66000      india      in
27  52  5.4      1935        False      44000      99000      india      in
28  39  3.4      1932        False      39000      88000      australia  aus
```

Coding Playground

```
In [ ] : import pandas as pd # We are not really calling pandas everytime
import numpy as np # Save case for numpy

# Splitting based on different separators for different locations
def split_location(row):
    location = row["Location"]
    if ',' in location:
        parts = location.split(',')
        country = parts[0].strip()
        initials = parts[1].strip()
    elif '-' in location:
        parts = location.split('-')
        country = parts[0].strip()
        initials = parts[1].strip()
    else:
        country = location.strip()
        initials = None # Or set initials as needed for cases with no separator
    return pd.Series({'Country': country, 'Initials': initials}) # Assigning each series with the values of p
articular row

url='https://raw.githubusercontent.com/ArjitDhall/PrepInsta-DA-Week-3/main/Assignment1.csv' # To make this cod
e available everywhere
df=pd.read_csv(url,encoding='unicode_escape') # To ensure proper format of data frame

df.drop_duplicates() # Dropping all the duplicate data if present

df["Easy Apply"] = df["Easy Apply"].str.replace('-1', 'FALSE') # Replacing -1 with False
mapping = {'TRUE': True, 'FALSE': False} # Converting strings to boolean values
df["Established"] = df["Established"].astype(str) # Convert Established to str data type
df = df.replace(-1, np.nan) # Replacing the remaining -1 with NaN Value

age_mean = df["Age"].mean() # For Replacing NaN values with mean of Age
age_mean = int(age_mean)
df["Age"] = df["Age"].replace(np.nan, age_mean)

rating_mean = df["Rating"].mean() # For Replacing NaN value with mean of Rating
rating_mean = int(rating_mean)
df["Rating"] = df["Rating"].replace(np.nan, rating_mean)

df["Established"] = df["Established"].replace('-1', 'Unknown')

df.drop_duplicates() # Drop remaining rows where Null value is present

df[["Starting Salary", "Ending Salary"]] = df["Salary"].str.split('-', 1, expand=True) # Splitting salary at '-'
df.drop(columns=["Salary"]) # No more need of Salary column

df["Starting Salary"] = df["Starting Salary"].str.strip("$") # Left Stripping '$'
df["Ending Salary"] = df["Ending Salary"].str.strip("$") # Mapped the string to respective
df["Starting Salary"] = df["Starting Salary"].str.replace('k', '000') # Replacing 'k' with '000'
df["Ending Salary"] = df["Ending Salary"].str.replace('k', '000')

df["Starting Salary"] = df["Starting Salary"].astype(int) # Convert Starting salary to int data type
df["Ending Salary"] = df["Ending Salary"].astype(int) # Convert Ending salary to int data type
df["Easy Apply"] = df["Easy Apply"].map(mapping).fillna(False).astype(bool) # Convert Easy Apply to int boolean ty
pe

# Apply the function to split 'Location' into 'Country' and 'Initials'
df[["Country", "Initials"]] = df.apply(split_location, axis=1) # Assigning the splitted values of Location to n
df.drop(columns=["Location"]) # No need of Location column anymore
df.reset_index(drop=True, inplace=True) # To restructure the index
df.drop(columns=["Index"], inplace=True)
df

<python>In[56]:AgeRatingEstablishedEasyApplyStartingSalaryEndingSalaryCountryInitials
except for the argument 'pat' will be keyword-only.
df[["Starting Salary", "Ending Salary"]] = df["Salary"].str.split('-', 1, expand=True) # Splitting salary at '-'

Out [ ] :
   age  rating  established  easy_apply  starting_salary  ending_salary  country  initials
0  44  5.4      1999        True      44000      99000      india      in
1  66  3.5      2002        True      55000      66000      new york   ny
2  39  4.0      unknown      False      77000      89000      new york   ny
3  64  4.4      1988        False      44000      99000      india      in
4  25  6.4      2002        False      44000      99000      australia  aus
5  44  1.4      1999        True      77000      89000      india      in
6  21  0.0      unknown      False      44000      99000      new york   ny
7  44  4.0      unknown      False      44000      99000      australia  aus
8  35  5.4      unknown      True      44000      99000      new york   ny
9  22  7.7      unknown      True      44000      99000      india      in
10  55  5.4      2008        True      10000      49000      india      in
11  44  6.7      2009        False      10000      49000      india      in
12  39  0.0      1999        False      44000      99000      india      in
13  25  4.0      2019        True      44000      99000      australia  aus
14  66  4.0      2020        True      44000      99000      australia  aus
15  44  3.0      1999        False      89000      101000     australia  aus
16  19  4.5      1984        False      19000      40000      india      in
17  39  5.3      1943        True      44000      99000      new york   ny
18  35  6.7      1954        True      44000      99000      new york   ny
19  32  3.3      1955        True      44000      99000      new york   ny
20  39  5.7      1944        True      44000      99000      new york   ny
21  35  5.0      1946        False      44000      99000      new york   ny
22  39  7.8      1988        True      55000      66000      new york   ny
23  39  2.4      1980        True      44000      99000      new york   ny
24  13  4.0      1987        False      44000      99000      new york   ny
25  55  0.0      1980        True      44000      99000      australia  aus
26  39  4.0      1934        True      55000      66000      india      in
27  52  5.4      1935        False      44000      99000      india      in
28  39  3.4      
```