Name: Sadneya Sadanand Samant

Roll no:46 Class:D15C

# **EXPERIMENT NO: 1**

**Aim**: Introduction to Data science and Data preparation using Pandas steps.

# Theory:

Data preparation is a crucial step in data science, involving cleaning and transforming raw data into an analyzable format. Using Pandas, we can perform operations such as handling missing values, encoding categorical data, and scaling numerical features. Proper preprocessing ensures the dataset is reliable for analysis and modeling by addressing inconsistencies, missing data, and outliers.

#### **Problem Statement:**

The Vehicle Safety Recall dataset, provided by NHTSA, contains 15 columns detailing various aspects of recall events, such as manufacturers, affected components, and corrective actions. This analysis focuses on:

- **Manufacturer Trends**: Identifying manufacturers prone to frequent recalls or specific defects.
- Impact Analysis: Understanding recall types affecting the largest populations and assessing average completion rates.
- Temporal Patterns: Detecting trends in recalls over time and seasonal spikes.
- **Safety Implications**: Investigating critical safety advisories like "Do Not Drive" or "Park Outside" and their resolution rates.

By cleaning the dataset and applying data preprocessing steps, the goal is to enhance its quality and draw actionable insights for stakeholders.

#### **Dataset Overview:**

The dataset provides detailed information about vehicle safety recalls managed by the National Highway Traffic Safety Administration (NHTSA). It contains 15 columns, each capturing specific aspects of recall events. Below is a breakdown of the columns and their relevance:

- **1. Report Received Date:** Date the recall was officially reported.
- 2. NHTSA ID: A unique identifier for each recall event.
- 3. Recall Link: A hyperlink to the recall details on the NHTSA website.
- **4. Manufacturer:** Name of the vehicle or product manufacturer responsible for the recall.

- **5. Subject:** Brief description of the recall issue.
- **6. Component:** The affected part of the vehicle/product (e.g., "POWER TRAIN").
- 7. Mfr Campaign Number: Manufacturer's internal reference for the recall.
- **8. Recall Type:** Type of product involved (e.g., vehicle, tire, or car seat).
- **9. Potentially Affected:** Number of units potentially impacted by the recall.
- **10. Recall Description:** Detailed explanation of the defect or issue.
- **11. Consequence Summary:** Description of the risks or consequences associated with the defect.
- **12. Corrective Action:** Steps taken to address the defect.
- **13. Park Outside Advisory:** Indicates whether there's an advisory to park outside for safety.
- **14. Do Not Drive Advisory:** Indicates whether there's an advisory not to drive the affected vehicle.
- **15. Completion Rate %:** Percentage of affected vehicles repaired or addressed.

## **Steps:**

## 1. Loading The Dataset

```
    [1] import pandas as pd

    [2] df = pd.read_csv('recalls.csv')
```

#### 2. Description of the dataset

#### a. Information about dataset

```
df.info()
→ ⟨class 'pandas.core.frame.DataFrame'>
    RangeIndex: 28671 entries, 0 to 28670
    Data columns (total 15 columns):
     # Column
                                                  Non-Null Count Dtype
                                                  28671 non-null object
     0 Report Received Date
     1 NHTSA ID
                                                  28671 non-null object
                                                  28671 non-null object
     2 Recall Link
                                                  28671 non-null object
28671 non-null object
     3 Manufacturer4 Subject
     5 Component
                                                 28671 non-null object
     6 Mfr Campaign Number
                                                 28624 non-null object
     7 Recall Type
                                                 28671 non-null object
     8 Potentially Affected
                                                 28630 non-null float64
                                                 26270 non-null object
23783 non-null object
        Recall Description
     10 Consequence Summary
     11 Corrective Action
                                                 26283 non-null object
     12 Park Outside Advisory
                                      28671 non-null object
                                                  28671 non-null object
     13 Do Not Drive Advisory
     14 Completion Rate % (Blank - Not Reported) 10007 non-null float64
    dtypes: float64(2), object(13)
    memory usage: 3.3+ MB
```

### **b.** Description of Dataset

```
# Get the dataset's shape and basic statistics
       print(f"Dataset Shape: {df.shape}")
       print(df.describe(include='all'))
 Dataset Shape: (28671, 15)
                                                                                                           Recall Link \
          Report Received Date
                                       NHTSA ID
                                                           count
                                                                                                                  28671
                             28671
                                           28671
                                                           unique
                                                                                                                  28671
 unique
                             10023
                                           28671
                                                           top
                                                                    Go to Recall (<a href="https://www.nhtsa.gov/recalls?nh">https://www.nhtsa.gov/recalls?nh</a>...
 top
                      10/17/2013
                                      25E002000
                                                           frea
 frea
                                42
                                                           mean
                                                                                                                    NaN
                               NaN
                                             NaN
 mean
                                                                                                                    NaN
 std
                               NaN
                                             NaN
                                                           min
 25%
                                NaN
                                             NaN
                                                           25%
                                                                                                                    NaN
 50%
                                NaN
                                             NaN
                                                           50%
                                                                                                                    NaN
 75%
                                NaN
                                             NaN
                                                           75%
                                                                                                                    NaN
 max
                                NaN
                                             NaN
                                                                                                                    NaN
                                                                                                           Recall Description
       Mfr Campaign Number Recall Type
                                       Potentially Affected
count
                    28624
                                28671
                                               2.863000e+04
                                                                 unique
unique
                    11341
                                                        NaN
                                                                         ON CERTAIN TRAILERS EQUIPPED WITH SEALCO SPRIN...
                                                                top
top
         NR (Not Reported)
                              Vehicle
                                                        NaN
                                                                freq
freq
                    16602
                                 24940
                                                        NaN
                                                                 mean
                                                                                                                           NaN
mean
                      NaN
                                  NaN
                                               4.572011e+04
                                                                                                                           NaN
std
                       NaN
                                  NaN
                                               3.730381e+05
                                                                min
                                                                                                                           NaN
min
                       NaN
                                  NaN
                                               0.000000e+00
                                                                25%
                                                                                                                           NaN
25%
                       NaN
                                  NaN
                                               9.900000e+01
                                                                50%
50%
                       NaN
                                  NaN
                                               6.860000e+02
                                                                                                                          NaN
75%
                       NaN
                                  NaN
                                               6.385500e+03
                                                                75%
                                                                                                                           NaN
                                               3.200000e+07
max
                                  NaN
                                      Consequence Summary
                                                                                                         Corrective Action \
count
                                                     23783
                                                                 count
unique
                                                                 unique
        RELEASE OF COOLANT UNDER CERTAIN CONDITIONS CO...
top
                                                                         DEALERS WILL EQUIP AIR SYSTEMS WITH A PRESSURE...
                                                                 top
freq
                                                       128
                                                                 freq
mean
                                                       NaN
                                                                                                                        NaN
                                                                 mean
std
                                                       NaN
                                                                 std
                                                                                                                        NaN
min
                                                       NaN
                                                                                                                        NaN
                                                                 min
25%
                                                       NaN
                                                                 25%
                                                                                                                        NaN
50%
                                                       NaN
                                                                 50%
                                                                                                                        NaN
75%
                                                       NaN
                                                                 75%
                                                                                                                        NaN
max
                                                       NaN
                                                                                                                        NaN
        Park Outside Advisory Do Not Drive Advisory
                                                                             Completion Rate % (Blank - Not Reported)
count
                                                                                                          10007.000000
                                                                    count
unique
                                2
                                                         2
                                                                    unique
top
                                                                    top
                                                                                                                    NaN
freq
                           28601
                                                     28510
                                                                    frea
                                                                                                                    NaN
mean
                              NaN
                                                       NaN
                                                                    mean
                                                                                                              67.874214
                                                                    std
                                                                                                              29.937993
std
                              NaN
                                                       NaN
                                                                    min
                                                                                                               0.000000
min
                              NaN
                                                       NaN
                                                                    25%
                                                                                                              48.350000
25%
                              NaN
                                                       NaN
                                                                    50%
                                                                                                              76.390000
50%
                              NaN
                                                       NaN
                                                                                                              93.765000
                                                                    75%
75%
                              NaN
                                                       NaN
                                                                                                             100.000000
max
                              NaN
```

# 3. Drop columns that aren't useful.

```
# Remove leading/trailing spaces from column names
df.columns = df.columns.str.strip()

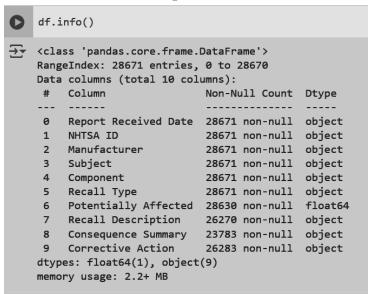
# List of columns to drop
cols = ["Recall Link", "Mfr Campaign Number", "Park Outside Advisory", "Do Not Drive Advisory", "Completion Rate % (Blank - Not Reported)"]
```

```
# Drop the columns that are present in the DataFrame
df = df.drop(cols, axis=1)

# Display the updated DataFrame
print(df.head())
```

```
Recall Description \
7
      Report Received Date
                                                          Manufacturer \
                                                                                       0 GKN Automotive (GKN) is recalling certain repl...
                                                                                       1 N&B Mobility Solutions LLC (Nivion) is recalli...
2 Forest River. Inc. (Forest 2:
                 01/14/2025 25E002000
                                                       GKN Automotive
                01/13/2025 25E001000 N&B Mobility Solutions LLC
01/13/2025 25V005000 Forest River, Inc.
01/13/2025 25V006000 Kia America, Inc.
                                                                                          Kia America, Inc. (Kia) is recalling certain 2...
                 01/13/2025 25V007000 Winnebago Industries, Inc.
                                                                                       4 Winnebago Industries, Inc. (Winnebago) is reca...
                                                                                                                         Consequence Summary \
                                                                                       0 A cracked or broken driveshaft can cause a los...
                                     Driveshaft Can Break
                                                                    POWER TRAIN
                                                                                       1 Inadequate clearance between DC busbars may ca...
          Charger Adapter May Cause Arcing or Shock Risk ELECTRICAL SYSTEM
                                                                                       2 A gas leak in the presence of an ignition sour...
                                                                       EQUIPMENT
      Cooktop Burner Tube May Crack and Cause Gas Leak
                                                                                       3 A loss of headlights and taillights can reduce...
           Loss of Headlights and Taillights/FMVSS 108 ELECTRICAL SYSTEM
                                                                                       4 A detached spare tire carrier can become a roa...
                           Spare Tire Carrier May Detach
                                                                      EOUIPMENT
                                                                                                                          Corrective Action
      Recall Type Potentially Affected \
                         130.0
                                                                                        0 GKN will reimburse the cost of a replacement d...
       Equipment
                                                                                       1 Nivion will replace the defective adapters, fr...
           Vehicle
Vehicle
                                     396.0
                                                                                          Owners are advised not to use the cooktop unti...
                                                                                       3 Dealers will update the BDC software, free of ...
                                  74469.0
                                                                                          Dealers will inspect, replace, and correctly t...
```

### Thus the columns now present in dataset are:

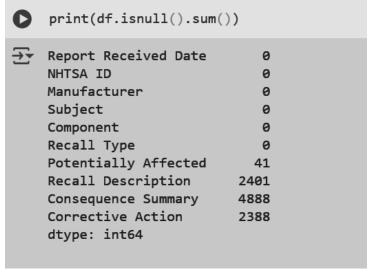


- 4. Take care of missing data.
  - a. Drop rows with maximum missing values.

```
print(f"Dataset Shape before Dropping Rows: {df.shape}")
# Drop rows with the highest number of missing values
threshold = len(df.columns) * 0.5 # Drop rows where over 50% of columns are missing
df = df.dropna(thresh=threshold)

print(f"Dataset Shape After Dropping Rows: {df.shape}")

Dataset Shape before Dropping Rows: (28671, 10)
Dataset Shape After Dropping Rows: (28671, 10)
```



# b. Handle Missing Data

Here above info says Potential Affected ,Recall Description ,Consequence Summary and corrective action contain some null values thus we need to handle missing data.

```
[12] # Fill missing numerical values with the median
     df['Potentially Affected'] = df['Potentially Affected'].fillna(df['Potentially Affected'].median())
     # Fill missing categorical values with a placeholder
     df['Recall Description'] = df['Recall Description'].fillna('Not Known')
     df['Consequence Summary'] = df['Consequence Summary'].fillna('Unknown')
     df['Corrective Action'] = df['Corrective Action'].fillna('Unknown')
     print(df.isnull().sum()) # Verify no missing values remain
→▼ Report Received Date
                             а
     NHTSA ID
     Manufacturer
    Subject
    Component
     Recall Type
     Potentially Affected 0
    Recall Description
Consequence Summary
                             0
     Corrective Action
     dtype: int64
```

#### 5. Create dummy variables

```
# Convert categorical columns into dummy variables
df = pd.get_dummies(df, columns=['Recall Type'], drop_first=True)
print(df.head())
```

```
Report Received Date NHTSA ID
                                               Manufacturer \
          01/14/2025 25E002000
                                             GKN Automotive
          01/13/2025 25E001000 N&B Mobility Solutions LLC
                                        Forest River, Inc.
          01/13/2025 25V005000
          01/13/2025 25V006000
                                         Kia America, Inc.
          01/13/2025 25V007000 Winnebago Industries, Inc.
                                          Subject
                                                          Component \
                            Driveshaft Can Break
                                                         POWER TRAIN
  Charger Adapter May Cause Arcing or Shock Risk ELECTRICAL SYSTEM
Cooktop Burner Tube May Crack and Cause Gas Leak
                                                           EOUIPMENT
      Loss of Headlights and Taillights/FMVSS 108 ELECTRICAL SYSTEM
                   Spare Tire Carrier May Detach
                                                           EOUIPMENT
 Potentially Affected
                                                      Recall Description
                 18.0 GKN Automotive (GKN) is recalling certain repl...
                130.0 N&B Mobility Solutions LLC (Nivion) is recalli...
                396.0 Forest River, Inc. (Forest River) is recalling...
              74469.0 Kia America, Inc. (Kia) is recalling certain 2...
                107.0 Winnebago Industries, Inc. (Winnebago) is reca...
                               Consequence Summary \
0 A cracked or broken driveshaft can cause a los...
 Inadequate clearance between DC busbars may ca...
 A gas leak in the presence of an ignition sour...
 A loss of headlights and taillights can reduce...
4 A detached spare tire carrier can become a roa...
                                Corrective Action Recall Type_Equipment \
0 GKN will reimburse the cost of a replacement d...
                                                                   True
1 Nivion will replace the defective adapters, fr...
                                                                   True
2 Owners are advised not to use the cooktop unti...
                                                                  False
3 Dealers will update the BDC software, free of ...
                                                                  False
4 Dealers will inspect, replace, and correctly t...
                                                                  False
   Recall Type_Tire Recall Type_Vehicle
             False
                                 False
             False
```

```
df.info()
<<class 'pandas.core.frame.DataFrame'>
    RangeIndex: 28671 entries, 0 to 28670
    Data columns (total 12 columns):
                               Non-Null Count Dtype
     # Column
         Report Received Date 28671 non-null object NHTSA ID 28671 non-null object
         Manufacturer
                              28671 non-null object
         Subject
                               28671 non-null object
                              28671 non-null object
         Component
         Potentially Affected 28671 non-null float64
         Recall Description
                               28671 non-null object
         Consequence Summary
                               28671 non-null object
     8 Corrective Action
                               28671 non-null object
         Recall Type_Equipment 28671 non-null bool
     10 Recall Type_Tire
                                28671 non-null
                                               bool
     11 Recall Type_Vehicle
                               28671 non-null bool
    dtypes: bool(3), float64(1), object(8)
    memory usage: 2.1+ MB
```

## 6. Find out outliers (manually)

```
# Specify the column to analyze for outliers
col = 'Potentially Affected'

# Calculate Q1, Q3, and IQR
Q1 = df[col].quantile(0.25)
Q3 = df[col].quantile(0.75)
IQR = Q3 - Q1

# Define lower and upper bounds
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Identify outliers
outliers = df[(df[col] < lower_bound) | (df[col] > upper_bound)]

# Display the outliers
print(f"Outliers in '{col}':")
print(outliers)
```

```
Outliers in 'Potentially Affected':
     Report Received Date NHTSA ID
                                                         Manufacturer \
               01/13/2025 25V006000
                                                    Kia America, Inc.
               01/06/2025 25V002000
                                                         Tesla, Inc.
14
               12/23/2024 24E110000
                                                        Horizon Global
21
               12/20/2024 24V957000
                                                    Ford Motor Company
22
               12/20/2024 24V954000
                                                    Ford Motor Company
               10/06/1966 66V004002
                                                   Ford Motor Company
               09/29/1966 66V003000 Honda (American Honda Motor Co.)
28668
               01/19/1966 66V032001
                                                  General Motors, LLC
28669
               01/19/1966 66V032003
                                                  General Motors, LLC
               01/19/1966 66V032004
28670
                                                  General Motors, LLC
                                               Subject \
            Loss of Headlights and Taillights/FMVSS 108
               Rearview Camera Image May Fail/FMVSS 111
      Tow Vehicle May Separate From Hitch Receiver Lock
                       High Pressure Fuel Pump May Fail
22
                 High Voltage Battery May Short Circuit
```

```
Loss of Headlights and Taillights/FMVSS 108
                Rearview Camera Image May Fail/FMVSS 111
       Tow Vehicle May Separate From Hitch Receiver Lock
                        High Pressure Fuel Pump May Fail
22
                  High Voltage Battery May Short Circuit
28658 INTERIOR SYSTEMS: RESTRAINT: BELT ANCHOR AND ATT...
                POWER TRAIN: TRANSMISSION: STANDARD: MANUAL
28666
                                          STEERING: COLUMN
28668
28669
                                          STEERING: COLUMN
                                          STEERING: COLUMN
28670
                  Component Potentially Affected \
          ELECTRICAL SYSTEM
                                           74469.0
       BACK OVER PREVENTION
                                          239382.0
14
            TRAILER HITCHES
                                          145431.0
       FUEL SYSTEM, DIESEL
ELECTRICAL SYSTEM
21
                                          295449.0
22
                                           20484.0
                 SEAT BELTS
                                           65000.0
                POWER TRAIN
28668
                   STEERING
                                          138878.0
28669
                   STEERING
                                           70644.0
28670
                   STEERING
                                           68184.0
```

```
Recall Description \
      Kia America, Inc. (Kia) is recalling certain 2...
      Tesla, Inc. (Tesla) is recalling certain 2024-...
      Horizon Global (Horizon) is recalling certain ...
      Ford Motor Company (Ford) is recalling certain...
21
22
      Ford Motor Company (Ford) is recalling certain...
28658
                                                Not Known
28666
                                                Not Known
28668
                                                Not Known
28669
                                                Not Known
28670
                                                Not Known
                                     Consequence Summary \
      A loss of headlights and taillights can reduce...
      A rearview camera that does not display an ima...
14
      A separated cap can allow the hitch to separat...
21
      High pressure Fuel pump failure can cause a lo...
      Battery failure can cause a loss of drive powe...
22
28658
                                                  Unknown
28666
                                                  Unknown
28668
                                                  Unknown
28669
                                                  Unknown
28670
                                                  Unknown
```

```
Corrective Action \
       Dealers will update the BDC software, free of ...
3
       Tesla released an over-the-air (OTA) software ...
7
       Dealers will replace the hitch receiver locks,...
       Dealers will update the powertrain control mod...
21
22
       Dealers will perform a battery energy control ...
28658
                                                    Unknown
28666
                                                    Unknown
28668
                                                    Unknown
28669
                                                    Unknown
28670
                                                    Unknown
       Recall Type_Equipment Recall Type_Tire
                                                   Recall Type_Vehicle
3
                        False
                                           False
                                                                   True
7
                        False
                                           False
                                                                   True
14
                         True
                                           False
                                                                  False
21
                        False
                                           False
                                                                   True
22
                                                                   True
                        False
                                           False
. . .
                                             . . .
                          . . .
                                                                    . . .
28658
                        False
                                           False
                                                                   True
28666
                                           False
                        False
                                                                   True
28668
                        False
                                           False
                                                                   True
28669
                        False
                                           False
                                                                   True
28670
                        False
                                           False
                                                                   True
[5063 rows x 12 columns]
```

#### 7. standardization and normalization of columns

```
from sklearn.preprocessing import StandardScaler, MinMaxScaler
    # Standardization: Transform data to have a mean of 0 and a standard deviation of 1
    standard_scaler = StandardScaler()
    df['Potentially Affected (Standardized)'] = standard_scaler.fit_transform(df[['Potentially Affected']])
    # Normalization: Scale data between 0 and 1
    min_max_scaler = MinMaxScaler()
    df['Potentially Affected (Normalized)'] = min_max_scaler.fit_transform(df[['Potentially Affected']])
    # Display the updated DataFrame
    print(df[['Potentially Affected', 'Potentially Affected (Standardized)', 'Potentially Affected (Normalized)']].head())
₹
       Potentially Affected Potentially Affected (Standardized) \
                      18.0
                      130.0
                                                      -0.122129
    1
                     396.0
                                                      -0.121415
                   74469.0
    3
                                                       0.077295
                     107.0
                                                       -0.122190
       Potentially Affected (Normalized)
                           5.625000e-07
    1
                           4.062500e-06
    2
                           1.237500e-05
    3
                            2.327156e-03
                            3.343750e-06
```

#### **Conclusion:**

This experiment demonstrated effective data cleaning and preparation techniques. Issues such as missing values, irrelevant data, and outliers were addressed, and the dataset was scaled for uniformity. These steps are essential for ensuring high-quality data and reliable model outcomes.