Data Cleaning is the process of identifying and fixing or removing errors and inconsistencies in data to improve its quality. This is a critical step in the data analysis pipeline to ensure that your analysis or model is built on accurate, complete, and reliable data. Data Cleaning is also called data cleansing or scrubbing -->Why is Data Cleaning Important? Improves accuracy of analysis and models
Ensures valid results from data insights Eliminates duplicate, corrupted, or incomplete data
 Helps maintain data consistency across sources

Deepika AIML Trainee | AI/ML bootcamp

**Data Cleaning** 

- --> Data Cleaning Tasks
- Removing Records: Deleting rows with missing values if they are relatively few and insignificant. • <u>Imputing Values:</u> Replacing missing values with estimated ones, such as the mean, median, or mode of the dataset.
- 2. Removing Duplicates:-Duplicates can skew analyses and lead to inaccurate results. Identifying and removing duplicate records ensures that each data point is unique and accurately represented.
- 3. Correcting Inaccuracies:-Data entry errors, such as typos or incorrect values, need to be identified and corrected. This can involve cross-referencing with other data sources or using validation rules to ensure data accuracy.

1. Handling Missing Data:-Missing data is a common problem in datasets. Strategies to handle missing data include:

- <u>Date formats</u> (2025-07-29 vs 29-07-2025) • Units (kg, lbs) • <u>Capitalization</u> (INDIA, India, india → India)
- -->Challenges in Data Cleaning
- Volume of Data: Large datasets can be challenging to clean due to their sheer size. Efficient techniques and tools are necessary to handle big data cleaning tasks.

4. Standardizing Data:-Standardization means bringing data to a common format, like:

- Complexity of Data: Data from diverse sources may have different structures and formats, making it difficult to clean and integrate.
- Continuous Process: Data cleaning is not a one-time task but an ongoing process. As new data is collected, it needs to be continually cleaned and maintained.

values

- The data cleaning process
- \*\* **1** Missing Irrelevant Duplication **Outliers**

data

PC 17758 108.9000 C105

7.2500

8.0500

22.3583

Fare Cabin Embarked

NaN

NaN

NaN

Q

Q

S

S

Q

С

S

Fare Cabin Embarked

NaN

S

С

0 SOTON/O.Q. 3101262

Ticket

330911

0 363272

240276

0 315154 8.6625

7538

0 330972 7.6292

Ticket

0 SOTON/O.Q. 3101262

1 3101298 12.2875 NaN

7.8292

7.0000

9.6875

9.2250

PC 17758 108.9000 C105

Ticket

392091

368364

59.4000

9.3500

7.7500

Fare Cabin Embarked

NaN

S

S

S

363272

240276

7.0000

315154 8.6625 NaN

9.2250

2657 7.2292 NaN

330972 7.6292 NaN

1 3101298 12.2875 NaN

248738 29.0000

0 A/4 48871 24.1500 NaN

17770 27.7208

248726 13.5000

Fare Cabin Embarked

С

С

S

S

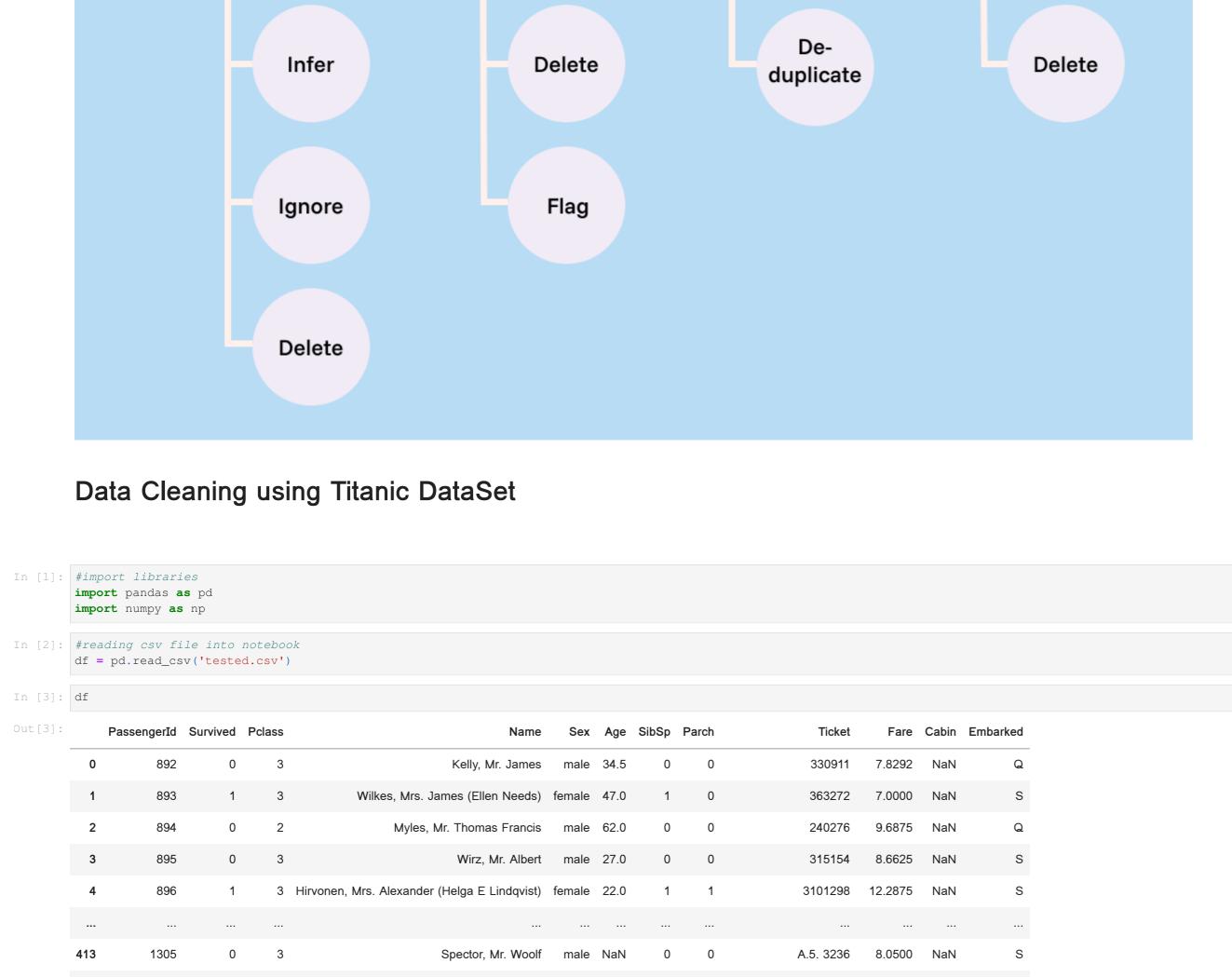
Q

7.2500

8.0500

22.3583 NaN

359309



## 418 rows × 12 columns In [4]: #print 1st to 7 columns df.head(7)

3

414

415

416

417

2

6

414

415

416

417

df.describe()

df.sample(5)

20

168

348

307

313

1306

1307

1308

1309

PassengerId Survived Pclass

894

896

898

1306

1307

1308

1309

PassengerId

PassengerId Survived Pclass

In [5]: #print last 4 column df.tail(4)

0

0

count 418.000000 418.000000 418.000000 332.000000 418.000000 418.000000 417.000000 1100.500000 0.363636 2.265550 30.272590 0.447368 0.392344 35.627188 120.810458 55.907576 0.481622 0.841838 14.181209 0.896760 0.981429 892.000000 0.000000 1.000000 0.170000 0.000000 0.000000 0.000000 996.250000 0.000000 1.000000 21.000000 0.000000 0.000000 7.895800 **50%** 1100.500000 0.000000 3.000000 27.000000 0.000000 0.000000 14.454200 3.000000 1.000000 0.000000 **75%** 1204.750000 1.000000 39.000000 31.500000 max 1309.000000 1.000000 3.000000 76.000000 8.000000 9.000000 512.329200

Rothschild, Mr. Martin

Aks, Master. Philip Frank

Wilkes, Mrs. James (Ellen Needs) female 47.0

Svensson, Mr. Johan Cervin male 14.0

Caldwell, Mr. Albert Francis male 26.0

Davies, Mr. John Samuel male 21.0

Connolly, Miss. Kate female 30.0

male 62.0

Wirz, Mr. Albert male 27.0

Myles, Mr. Thomas Francis

3 Hirvonen, Mrs. Alexander (Helga E Lindqvist) female 22.0

3 Abrahim, Mrs. Joseph (Sophie Halaut Easu) female 18.0

Carr, Miss. Jeannie female 37.00

1 Cassebeer, Mrs. Henry Arthur Jr (Eleanor Genev... female

SibSp

Oliva y Ocana, Dona. Fermina female 39.0

Ware, Mr. Frederick

Peter, Master. Michael J

Kelly, Mr. James

Svensson, Mr. Johan Cervin male 14.0

Connolly, Miss. Kate female 30.0

male 38.5

male NaN

Wilkes, Mrs. James (Ellen Needs)

3 Hirvonen, Mrs. Alexander (Helga E Lindqvist) female 22.0

1 Oliva y Ocana, Dona. Fermina female 39.0

Ware, Mr. Frederick

Peter, Master. Michael J

Saether, Mr. Simon Sivertsen

In [6]: # get a summary of the DataFrame, statistics of the numerical columns

Survived

In [7]: # get a random sample of 5 rows from the DataFrame

PassengerId Survived Pclass

1060

1240

1199

1205

897

898

899

901

df.isnull().sum()

Out[10]: PassengerId

Age

Survived Pclass Name

0

0

0

In [10]: # checking for missing values in the DataFrame

Myles, Mr. Thomas Francis

male 38.5

male NaN

male NaN

male 62.0

Sex Age SibSp Parch

Parch

Fare

Sex Age SibSp Parch

male 55.00

male 24.00

Wirz, Mr. Albert male 27.0

Sex Age SibSp Parch

Saether, Mr. Simon Sivertsen

In [8]: # check the data types and null values in the DataFrame df.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 418 entries, 0 to 417 Data columns (total 12 columns): Non-Null Count Dtype Column PassengerId 418 non-null int64 Survived 418 non-null 418 non-null Pclass 418 non-null Name 418 non-null Age 332 non-null 418 non-null SibSp Parch 418 non-null Ticket 418 non-null Fare 417 non-null 10 Cabin 91 non-null 11 Embarked 418 non-null dtypes: float64(2), int64(5), object(5) memory usage: 39.3+ KB In [9]: #slicing df.iloc[1:10] PassengerId Survived Pclass Sex Age SibSp Parch

SibSp Parch Ticket Fare Cabin Embarked dtype: int64 In [11]: # particaluerly checking for missing values in the 'Ticket' columns df['Ticket'].isnull().sum() Out[11]: np.int64(0) In [12]: # multiple columns check for missing values df[['Ticket', 'Fare', 'Age']].isnull().sum() Out[12]: Ticket Fare dtype: int64 In [13]: # check for duplicate rows in the DataFrame df.duplicated().sum() Out[13]: np.int64(0) In [14]: #check shape df.shape Out[14]: (418, 12) In [15]: df['Name'].head() Out[15]: 0 Kelly, Mr. James Wilkes, Mrs. James (Ellen Needs)

• Listwise deletion:- removes entire rows from the dataset if any single value is missing or erroneous in that row.

Sex Age SibSp Parch

male 27.0

male NaN

male NaN

The behavior will change in pandas 3.0. This implace method will never work because the intermediate object on which we are setting values always behaves as a copy.

Sex Age SibSp Parch

male 62.0

male 27.0

female 22.0

male 38.5

male NaN

male NaN

0

0

Kelly, Mr. James

Wirz, Mr. Albert

Wilkes, Mrs. James (Ellen Needs)

3 Hirvonen, Mrs. Alexander (Helga E Lindqvist)

Myles, Mr. Thomas Francis

Kelly, Mr. James

Wilkes, Mrs. James (Ellen Needs) female 47.0

Wirz, Mr. Albert

Spector, Mr. Woolf

Peter, Master. Michael J

Myles, Mr. Thomas Francis

3 Hirvonen, Mrs. Alexander (Helga E Lindqvist) female 22.0

Fare Cabin Embarked

NaN

NaN

NaN

NaN

Fare Cabin Embarked

0

0

0

Q

S

S

С

Ticket

330911

363272

240276

315154

3101298

A.5. 3236

359309

**END** 

THANK YOU!

2668

SOTON/O.Q. 3101262

7.8292

7.0000

9.6875

8.6625

12.2875

8.0500

7.2500

8.0500

22.3583

PC 17758 108.9000 C105

Q

Q

S

S

С

330911

363272

240276

A.5. 3236

SOTON/O.Q. 3101262

7.8292

9.6875

7.0000 NaN

8.6625 NaN

12.2875 NaN

8.0500 NaN

PC 17758 108.9000 C105

7.2500

8.0500

22.3583 NaN

• Pairwise Deletion:-Pairwise deletion keeps all available data for each analysis, only excluding missing values when necessary.

Myles, Mr. Thomas Francis

21. , nan, 46. , 23. , 63. , 24. , 35. , 45. , 55. , 9. , 48. , 50. , 22.5 , 41. , 33. , 18.5 , 25. , 39. 60. , 36. , 20. , 28. , 10. , 17. , 32. , 13. , 31. , 29. , 28.5 , 32.5 , 6. , 67. , 49. , 2. , 76. , 43. , 16. , 1. , 12. , 42. , 53. , 26.5 , 40. , 61. , 60.5 , 7. , 15. , 54. , 64. , 37. , 34. , 11.5 , 8. , 0.33, 38. , 57. , 40.5 , 0.92, 19. , 36.5 , 0.75, 0.83, 58. , 0.17, 59. , 14.5 , 44. , 5. , 51. , 3. , 38.5 ])

4 Hirvonen, Mrs. Alexander (Helga E Lindqvist)

Out[16]: array([34.5 , 47. , 62. , 27. , 22. , 14. , 30. , 26. , 18. ,

Handle these errors or anomalies

In [18]: | # drop rows with missing values in 'Ticket' and 'Age' columns df\_pairwise\_deletion = df.dropna(subset=['Fare'])

Name: Name, dtype: object

In [17]: # drop rows with any missing values

In [19]: #df.isnull().sum()

Survived Pclass Name Sex

Out[19]: PassengerId

SibSp Parch Ticket Fare Cabin

Embarked dtype: int64

Out[20]: PassengerId 0 Survived Pclass

> Sex Age SibSp Parch Ticket Fare Cabin Embarked dtype: int64

In [21]: df.isnull().sum()

Survived Pclass

Embarked dtype: int64

2

4

413

417

Out[25]: 0

PassengerId Survived Pclass

892

896

1305

1309

In [23]: # calculate the mean of the 'Quantity' column mean\_quantity=df['Survived'].mean()

df1['Fare'].fillna(mean\_quantity)

In [27]: # round the mean value to 0 decimal places

PassengerId Survived Pclass

892

893

894

895

896

418 rows × 12 columns

median\_customer\_id

round(mean\_quantity)

7.8292 7.0000 9.6875 8.6625 12.2875

In [25]: # fill missing values in 'Quantity' with the mean value

418 rows × 12 columns

Out[23]: np.float64(0.36363636363636363)

mean\_quantity

df1 = df.copy()

In [26]: mean\_quantity//1

Out [26]: np.float64(0.0)

In [29]: df1

2

3

4

Out[29]:

In [22]: df

Out[22]:

Out[21]: PassengerId

Sex

df\_listwise\_deletion = df.dropna()

df\_pairwise\_deletion.isnull().sum()

326

In [20]: df\_listwise\_deletion.isnull().sum()

In [16]: df['Age'].unique()

Wirz, Mr. Albert

## Age SibSp Parch Ticket Fare 327 Cabin

**IMPUTATION** - impute - filling the values

414 1306 Oliva y Ocana, Dona. Fermina female 39.0 415 1307 Saether, Mr. Simon Sivertsen male 38.5 416 1308 Ware, Mr. Frederick male NaN

3

In [24]: # create a copy of the DataFrame to avoid modifying the original DataFrame

413 8.0500 414 108.9000 415 7.2500 416 8.0500 22.3583 Name: Fare, Length: 418, dtype: float64

In [28]: # fill missing values in 'Quantity' with the rounded mean value df1['Cabin'].fillna(round(mean\_quantity), inplace=True) C:\Users\Deepika\AppData\Local\Temp\ipykernel\_15300\1054803618.py:2: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method(value) instead, to perform the operation inplace on the original object. df1['Cabin'].fillna(round(mean\_quantity), inplace=True)

2

413 1305 Spector, Mr. Woolf Oliva y Ocana, Dona. Fermina 414 1306 415 1307 Saether, Mr. Simon Sivertsen 1308 416 Ware, Mr. Frederick 417 1309 3 Peter, Master. Michael J

In [30]: # calculate the median of the 'PassengerId' column

median\_customer\_id = df['PassengerId'].median()

0

Out[30]: np.float64(1100.5) In [31]: # calculate the mode of the 'PassengerId' column mode\_customer\_id = df['PassengerId'].mode()[0] # [] used here to get the first value of the mode series # [1] next mode Out[31]: np.int64(892)