



# ADNOC Accelerator Programme Artificial Intelligence Сонокт 2

# **Data Transformation and Cleaning**

© 2025 World Wide Technology, Inc. All rights reserved.

# **Data Transformation and Cleaning**





# Data processing is a crucial step in the Machine Learning lifecycle





# Data can be quantitative or qualitative, and have further subtypes



#### Datasets are explored through basic statistics and visualisation

#### Data columns

	Pipeline Location	Pipeline Type	Cause Category	Cause Subcategory	Net Loss (Barrels)	All Costs	
	ONSHORE	UNDERGROUND	CORROSION	INTERNAL	0.00	5065	
	ONSHORE	UNDERGROUND	MATERIAL/WELD/EQUIP FAILURE	CONSTRUCTION, INSTALLATION OR FABRICATION-RELATED	0.00	216121	
	ONSHORE	TANK	ALL OTHER CAUSES	MISCELLANEOUS	0.00	16200	
	ONSHORE	ABOVEGROUND	MATERIAL/WELD/EQUIP FAILURE	MALFUNCTION OF CONTROL/RELIEF EQUIPMENT	0.00	32477	
Missing	ONSHORE	UNDERGROUND	INCORRECT OPERATION	DAMAGE BY OPERATOR OR OPERATOR'S CONTRACTOR	25.00	84600	
wissing	ONSHORE	UNDERGROUND	EXCAVATION DAMAGE	THIRD PARTY EXCAVATION DAMAGE	0.00	709351	
values	ONSHORE	ABOVEGROUND	INCORRECT OPERATION	INCORRECT VALVE POSITION	0.00	16169	
	ONSHORE	ALGROUND	CORROSION	EXTERNAL	0.00	16025	
			CORROSION	EXTERNAL	0.10	1450000	
	ONSHOP	NaN	RIAL/WELD/EQUIP FAILURE	CONSTRUCTION, INSTALLATION OR FABRICATION-RELATED	0.00	40000	
	ONSHORE	ERGROUND	CORROSION	INTERNAL	0.00	17485	
	ONSHORE	ABON	CORROSION	INTERNAL	0.36	5040	
Categorical	ONSHORE	ABOVEGROUND	INCORRECT OPERATION	OVERFILL/OVERFLOW OF TANK/VESSEL/SUMP	0.00	11607	
columns	ONSHORE	ABOVEGROUND	MATERIAL/WELD/EQUIP FAILURE	MALFUNCTION OF CONTROL/RELIEF EQUIPMENT	0.35	6350	
••••••	ONSHORE	ABOVEGROUND	MATERIAL/WELD/EQUIP FAILURE	PUMP OR PUMP-RELATED EQUIPMENT	0.00	3770	
	ONSHORE	ABOVEGROUND	OTHER OUTSIDE FORCE DAMAGE	VEHICLE NOT ENGAGED IN EXCAVATION	0.10	150001	
		UNDERGROUND	EXCAVATION DAMAGE	OPERATOR/CONTRACTOR EXCAVATION DAMAGE	0.00	40750	
	ONSHORE	ABOVEGROUND	OTHER OUTSIDE FORCE DAMAGE	VEHICLE NOT ENGAGED IN EXCAVATION	976.00	776753	
	ONSHORE	ABOVEGROUND	MATERIAL/WELD/EQUIP FAILURE	NON-THREADED CONNECTION FAILURE	3.00	1270	
	ONSHORE	ABOVEGROUND	INCORRECT OPERATION	INCORRECT INSTALLATION	0.35	2034	

Numerical columns

#### Statistics can be summarised in different ways





#### **Provide Different Summary Stats**

<pre>df_sample.describe()</pre>		
	Net Loss (Barrels)	All Costs
count	2795.000000	2.795000e+03
mean	132.194050	8.340332e+05
std	1185.019252	1.657830e+07
min	0.000000	0.000000e+00
25%	0.000000	5.039500e+03
<b>50%</b>	0.000000	2.312900e+04
75%	2.000000	1.172325e+05
max	30565.000000	8.405261e+08

© 2025 World Wide Technology, Inc. All rights reserved.

з **Ж** 

© 2025 World Wide Technology, Inc. All rights reserved

• X

#### **Count unique categories**

#### df\_sample.nunique()

Pipeline Location	2
Pipeline Type	4
Cause Category	7
Cause Subcategory	38
Net Loss (Barrels)	443
All Costs	2279

#### **Check missing values**

df\_sample.isna().sum()

Pipeline Location	0
Pipeline Type	18
Cause Category	0
Cause Subcategory	0
Net Loss (Barrels)	0
All Costs	0



#### Show frequency of unique values

df\_sample[["Cause Category"]].value\_counts()

Cause Category			
MATERIAL/WELD/EQUIP FAILURE	1435		
CORROSION	592		
INCORRECT OPERATION	378		
ALL OTHER CAUSES 118			
NATURAL FORCE DAMAGE	118		
EXCAVATION DAMAGE 97			
OTHER OUTSIDE FORCE DAMAGE	57		

#### **Provide Different Summary Stats**

df\_sample.describe()

	Net Loss (Barrels)	All Costs
count	2795.000000	2.795000e+03
mean	132.194050	8.340332e+05
std	1185.019252	1.657830e+07
min	0.000000	0.000000e+00
25%	0.000000	5.039500e+03
<b>50%</b>	0.000000	2.312900e+04
75%	2.000000	1.172325e+05
max	30565.000000	8.405261e+08

# Summary Statistics provide a quick overview of your dataset



# **Central Tendency shows the typical or central value of dataset**



Dataset on house prices

House Number	House Price (AED)
101	1,200,000
102	1,100,000
103	1,250,000
104	1,100,000
105	1,200,000
106	1,300,000
107	1,150,000
108	1,300,000
109	1,400,000
110	3,500,000

 If you had to report the typical house price from this dataset, what value would you choose?

You can use a measure of central tendency – mean, median, or mode



House Number	House Price (AED)
101	1,200,000
102	1,100,000
103	1,250,000
104	1,200,000
105	1,200,000
106	1,300,000
107	1,150,000
108	1,300,000
109	1,400,000
110	3,500,000





House Number	House Price (AED)
101	1,200,000
102	1,100,000
103	1,250,000
104	1,200,000
105	1,200,000
106	1,300,000
107	1,150,000
108	1,300,000
109	1,400,000
110	3,500,000





House Number	House Price (AED)
102	1,100,000
107	1,150,000
101	1,200,000
104	1,200,000
105	1,200,000
103	1,250,000
106	1,300,000
108	1,300,000
109	1,400,000
110	3,500,000

#### Median

#### First, arrange in order

Then, find the middlemost value (1,200,000+1,250,000)/2 = 1,225,000



House Number	House Price (AED)
101	1,200,000
102	1,100,000
103	1,250,000
104	1,200,000
105	1,200,000
106	1,300,000
107	1,150,000
108	1,300,000
109	1,400,000
110	3,500,000



Test your knowledge!

# Which would help you identify the month with the most holidays?

- A. Mean
- B. Median
- C. Mode

Test your knowledge!

# Which would help you identify the month with the most holidays?

A. Mean

# B. Median

# C. Mode





House Number	House Price (AED)
101	1,200,000
102	1,100,000
103	1,250,000
104	1,200,000
105	1,200,000
106	1,300,000
107	1,150,000
108	1,300,000
109	1,400,000
110	3,500,000

#### Range

Difference between largest and smallest values

 $Range(x) = \max(x) - \min(x)$ 

Here, range would be

(3,500,000-1,100,000) = 2,400,000



Simply shows the full spread between smallest and largest values



House Number	House Price (AED)
102	1,100,000
107	1,150,000
101	1,200,000
104	1,200,000
105	1,200,000
103	1,250,000
106	1,300,000
108	1,300,000
109	1,400,000
110	3,500,000

#### Inter-quartile Range





Shows the spread of middle 50% of data while ignoring **outliers** 



House Number	House Price (AED)
102	1,100,000
107	1,150,000
101	1,200,000
104	1,200,000
105	1,200,000
103	1,250,000
106	1,300,000
108	1,300,000
109	1,400,000
110	3,500,000

#### Inter-quartile Range

Here, inter-quartile range would be

IQR=Q3-Q1=1,300,000-1,200,000=100,000

#### Shows the spread of middle 50% of data while ignoring **outliers**



Standard Deviation = Square Root of  $\sum_{i=1}^{N} (X_i - Mean)^2$ 



House Number	House Price (AED)	
101	1320000	
102	1230000	
103	1350000	
104	1480000	
105	1210000	
106	1210000	
107	1490000	
108	1370000	
109	1180000	
110	1330000	
Std Dev	104790	



House Number	House Price (AED)	
101	1000000	
102	1000000	
103	1400000	
104	500000	
105	500000	
106	900000	
107	700000	
108	1500000	
109	700000	
110	500000	
Std Dev	343656	

Mean is 1275000 for both



Shows how spread-out data points are around the average, here, 1275000

# Missing values can bias results and reduce model accuracy



#### Checking for missing values



# Missing numerical values can be filled with numerical averages





# Missing categorical values can be filled with mode or placeholder



#### Handling Outliers is essential to prepare your dataset

#### **A** What is an Outlier?

An outlier is a data point that deviates *significantly* from the rest of the dataset.

House Number	House Price (AED)
101	1,200,000
102	1,100,000
103	1,250,000
104	1,200,000
105	1,200,000
106	1,300,000
107	1,150,000
108	1,300,000
109	1,400,000
110	3,500,000

Remember this house that was pulling up your mean? That's your outlier!

#### Handling Outliers is essential to prepare your dataset

#### **A** What is an Outlier?

An outlier is a data point that deviates *significantly* from the rest of the dataset.

House Number	House Price (AED)
101	1,200,000
102	1,100,000
103	1,250,000
104	1,200,000
105	1,200,000
106	1,300,000
107	1,150,000
108	1,300,000
109	1,400,000
110	3,500,000

How do we decide if the deviation is significant?



#### **Outliers are those values which lie outside the inter-quartile range**





# Outliers can be rectified easily with the help of Python





#### Outliers can be rectified easily with the help of Python



29

#### Data can be rescaled for better outlier handling



# Scales values between 0 and 1 $x' = \frac{x - \min(x)}{\max(x) - \min(x)}$

#### Data can be rescaled for better outlier handling







Best for data within a range, free from outliers, and preserving relative distances



#### Data can be normalised to scale for better outlier handling



House Number	House Price (AED)
101	1,200,000
102	1,100,000
103	1,250,000
104	1,200,000
105	1,200,000
106	1,300,000
107	1,150,000
108	1,300,000
109	1,400,000
110	3,500,000



#### Data can be normalised to scale for better outlier handling



House Number	House Price (Standardised)
101	-0.380
102	-0.526
103	-0.307
104	-0.380
105	-0.380
106	-0.234
107	-0.453
108	-0.234
109	-0.088
110	1





This method is used often, especially for data that assume normal distribution

#### Feature Engineering can make input data suitable for ML

Input Data to a Machine Learning Model



#### **Ordinal Data can be treated with Label Encoding**

#### Label Encoding (Integer Encoding) Replaces each unique category with an integer

Index	Pipeline Type
0	ABOVEGROUND
1	UNDERGROUND
2	ABOVEGROUND
3	ABOVEGROUND
4	TANK



Index	Pipeline Type Encoded
0	0
1	1
2	0
3	0
4	2

#### Nominal categories are best handled with One-Hot Encoding

#### One-Hot Encoding (OHE) Creates a binary column for each category

Index	Pipeline Type
0	ABOVEGROUND
1	UNDERGROUND
2	ABOVEGROUND
3	ABOVEGROUND
4	UNDERGROUND

Index	Pipeline Type Aboveground	Pipeline Type Underground
0	1	0
1	0	1
2	1	0
3	1	0
4	0	1

#### **Correlation issues can destabilise and slow down your model**

Correlation is simply the relationship between two variables





#### Some features you chose may be related to each other

All Costs	Net Loss (Barrels)	Cause Subcategory	Cause Category	Pipeline Type	<b>Pipeline Location</b>
5065	0.00	INTERNAL	CORROSION	UNDERGROUND	ONSHORE
216121	0.00	CONSTRUCTION, INSTALLATION OR FABRICATION-RELATED	MATERIAL/WELD/EQUIP FAILURE	UNDERGROUND	ONSHORE
16200	0.00	MISCELLANEOUS	ALL OTHER CAUSES	TANK	ONSHORE
32477	0.00	MALFUNCTION OF CONTROL/RELIEF EQUIPMENT	MATERIAL/WELD/EQUIP FAILURE	ABOVEGROUND	ONSHORE
84600	25.00	DAMAGE BY OPERATOR OR OPERATOR'S CONTRACTOR	INCORRECT OPERATION	UNDERGROUND	ONSHORE
709351	0.00	THIRD PARTY EXCAVATION DAMAGE	EXCAVATION DAMAGE	UNDERGROUND	ONSHORE
16169	0.00	INCORRECT VALVE POSITION	INCORRECT OPERATION	ABOVEGROUND	ONSHORE
16025	0.00	EXTERNAL	CORROSION	ABOVEGROUND	ONSHORE
1450000	0.10	EXTERNAL	CORROSION	NaN	OFFSHORE
40000	0.00	CONSTRUCTION, INSTALLATION OR FABRICATION-RELATED	MATERIAL/WELD/EQUIP FAILURE	UNDERGROUND	ONSHORE
17485	0.00	INTERNAL	CORROSION	UNDERGROUND	ONSHORE
5040	0.36	INTERNAL	CORROSION	ABOVEGROUND	ONSHORE
11607	0.00	OVERFILL/OVERFLOW OF TANK/VESSEL/SUMP	INCORRECT OPERATION	ABOVEGROUND	ONSHORE
6350	0.35	MALFUNCTION OF CONTROL/RELIEF EQUIPMENT	MATERIAL/WELD/EQUIP FAILURE	ABOVEGROUND	ONSHORE
3770	0.00	PUMP OR PUMP-RELATED EQUIPMENT	MATERIAL/WELD/EQUIP FAILURE	ABOVEGROUND	ONSHORE
150001	0.10	VEHICLE NOT ENGAGED IN EXCAVATION	OTHER OUTSIDE FORCE DAMAGE	ABOVEGROUND	ONSHORE
40750	0.00	OPERATOR/CONTRACTOR EXCAVATION DAMAGE	EXCAVATION DAMAGE	UNDERGROUND	ONSHORE
776753	976.00	VEHICLE NOT ENGAGED IN EXCAVATION	OTHER OUTSIDE FORCE DAMAGE	ABOVEGROUND	ONSHORE
1270	3.00	NON-THREADED CONNECTION FAILURE	MATERIAL/WELD/EQUIP FAILURE	ABOVEGROUND	ONSHORE
2034	0.35	INCORRECT INSTALLATION	INCORRECT OPERATION	ABOVEGROUND	ONSHORE

Remember this dataset we saw?

#### Handling correlation is necessary for your model to be reliable

Let's say you chose two features to incorporate in you model: Unintentional Release and Net Loss



#### Handling correlation is necessary for your model to be reliable

# These two features 'Unintentional Release' and 'Net Loss' show a very high correlation



ā 0.



#### Handling correlation is necessary for your model to be reliable

Try to keep only one of the correlated features!



Redundant information

Slower Learning

Misleading Results

# **Data Transformation and Cleaning**

In this session, we covered:





Identifying and understanding data types

Using summary statistics to describe a dataset

Cleaning data by handling missing values and outliers



Transforming data using normalisation and standardisation



Using feature engineering to improve model training