

Exploratory Data Analysis

Questions to always ask before starting	<p>Q1. How big is your data? : <code>df.shape</code></p> <p>Q2. How does the data look like: <code>df.head()</code> , <code>df.sample(5)</code> → will display random rows</p> <p>Q3. What is the data type of cols : <code>df.info()</code> → fix any dtype issues and improve memory.</p> <p>Q4. Are there any missing values: <code>df.isnull().sum()</code> → gives a list of total null values in every column</p> <p>Q5. How does the data look like mathematically? ya phir descriptive statistics? Ans: <code>df.describe()</code> → only useful for numerical values.</p> <p>Q6. Are there duplicate values? Ans: <code>df.duplicated().sum()</code></p> <p>Q7. What is the correlation b/w cols? <code>df.corr()['column-name']</code> → gives a relation with every column</p>
univariate analysis	<p>technique that focuses on examining and summarising a single variable in isolation.</p> <div><div><p>categorical data</p><ol style="list-style-type: none">1. bar chart2. pie chart</div><div><p>numerical data</p><p>histogram</p><p>box plot</p><p>line chart</p><p>distplot</p></div></div>
bivariate analysis	<p>statistical analysis of 2 variables to determine the relationship between them. It involves examining the distribution, association and interaction b/w them.</p>
pandas profiling	<p>python library for EDA.</p>
profile report	<div><pre>from pandas_profiling import ProfileReport prof = ProfileReport(df) prof.to_file(output_file = 'output.html')</pre></div> <p>→ generates a web page for the dataset analysis. has overview, warnings and much more at a glance instantly.</p>

Feature Engineering > Feature Scaling

definition	is the process of transforming raw data into a format that is suitable for ML algorithms to extract meaningful patterns and classifications.			
overview	<div>Transformation</div> <ul style="list-style-type: none"> ① missing values ② categorical features ③ outlier detection ④ Feature Scaling 	<div>Construction</div> <ul style="list-style-type: none"> ↳ create a new feature from existing ones. 	<div>Selection</div> <ul style="list-style-type: none"> ↳ for optimising model 	<div>Extraction</div>
feature scaling ↳ data normalisation	the process of transforming the values of different features (variables) in a dataset to a common scale. This is done to ensure that all features have a similar magnitude and distribution.			
when to use feature scaling?	<ol style="list-style-type: none"> ML algorithms that rely on distance metrics like KNN or SVM. algorithms that use gradient based optimisation such as regression as it helps in achieving faster convergence and prevents certain features from dominating the optimization process due to their larger magnitude. 			
tip	check the documentation or specific guidelines for the ML algorithm being used to determine if feature scaling is necessary.			
types	standardisation, normalisation			
Standardization ↳ Z-score normalisation	$X'_i = \frac{X_i - \bar{X}}{\sigma}$ is a feature scaling technique that transforms the values of features to have a mean of 0 and a standard deviation of 1.			
when to use	distance based algorithm, linear models, PCA and regularisation			
when not to use	decision trees, random forest and sparse data			
geometric intuition	data is centered around the mean			
code	<ol style="list-style-type: none"> Train test split <pre>from sklearn.preprocessing import StandardScaler scaler = StandardScaler() scaler.fit(X_train) X_scaled_train = scaler.transform(X_train) X_scaled_test = scaler.transform(X_test)</pre> 			
impact of outliers	<p>outliers affect both mean and standard deviation causing the standardisation data to be skewed.</p> <p>Alternative methods like robust scaling or normalisation are better if your data has a lot of outliers.</p>			

X_{train} df
↓ transformed
↓
 X_{train} scaled is np array.

} transformation is on entire dataset.

Feature Scaling > Normalisation

minmax scaling	rescales the values of features to a specific range. $X_{norm} = (X_i - X_{min}) / (X_{max} - X_{min})$ it does not change the distribution or shape the data. It only rescales. range will always be in 0-1.
mean normalisation	$X_i' = \frac{X_i - X_{mean}}{X_{max} - X_{min}}$ mean centering.
max absolute scaling ↳ use when data is sparse	also known as max-min scaling, feature scaling technique that rescales the values of a feature to a range between -1 and 1. $X_{scaled} = X / \max_abs$. 'max_abs' is the maximum absolute value. It does not center the data around zero or adjust the variance.
median-mad scaling ↳ robust scaling	technique that rescales the values of features based on their robust statistics, making it less sensitive to outliers compared to other scaling methods. $X_{scaled} = (x - M) / MAD$ where M is median and MAD = M absolute deviation used when outliers are many.
how is it different from standardisation	most problems will be solved by standardisation. Normalisation tends to compress or stretches the original distribution. It is also more sensitive to outliers. Use it in image processes.

Encoding Categorical data

definition	refers to the process of converting categorical variables into a numerical representation.															
types	one hot encoding (dummy coding), Label Encoding, Ordinal Encoding and Binary Encoding.															
label encoding ↳ used for target value	Each unique category is assigned to different integer value. For example, consider a categorical variable "Color" with categories 'Red', 'Green' can be assigned 0, 1 and so on for other colors.															
ordinal encoding	similar to label encoding however ordinal encoding takes into account the inherent order or hierarchy present in the categories. example: PhD > Masters then the label assigned would be 3, 2. from sklearn.preprocessing import OrdinalEncoder oe = OrdinalEncoder(categories=['Poor', 'Average', 'Good']) oe.fit(X_train) X_train = oe.transform(X_train) X_test = oe.transform(X_test)															
one hot encoding	a technique used to convert categorical variables into binary vectors. It creates a binary columns for each unique category within a categorical feature. this may introduce high dimensionality to the dataset, especially if there are many unique categories within a variable.															
multicollinearity or dummy variable trap.	2 or more predictor variables in a regression model are highly correlated <table border="1"><thead><tr><th>Y</th><th>B</th><th>R</th></tr></thead><tbody><tr><td>1</td><td>0</td><td>0</td></tr><tr><td>0</td><td>1</td><td>0</td></tr><tr><td>0</td><td>0</td><td>1</td></tr><tr><td>1</td><td>0</td><td>0</td></tr></tbody></table> ↳ there is a mathematical relationship formed i.e. $\sum YBR = 1$. ↳ this is multicollinearity <u>AVOID</u> → to resolve this → if n dimensions are created keep n-1. Delete the 1st column.	Y	B	R	1	0	0	0	1	0	0	0	1	1	0	0
Y	B	R														
1	0	0														
0	1	0														
0	0	1														
1	0	0														
dimensionality																