

Data Collection and Preprocessing Phase

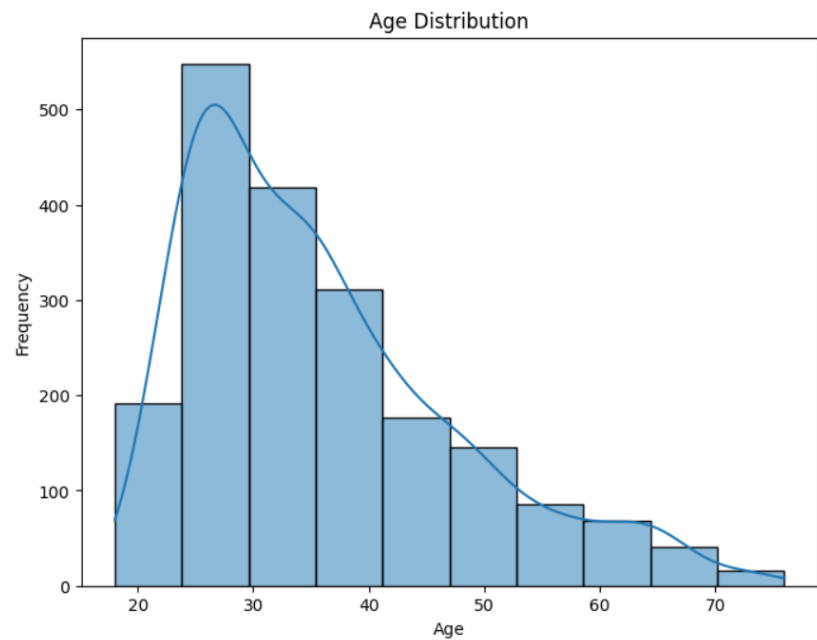
Date	15 March 2024
Team ID	LTVIP2024TMID25001
Project Title	Customer Segmentation Using Machine Learning
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

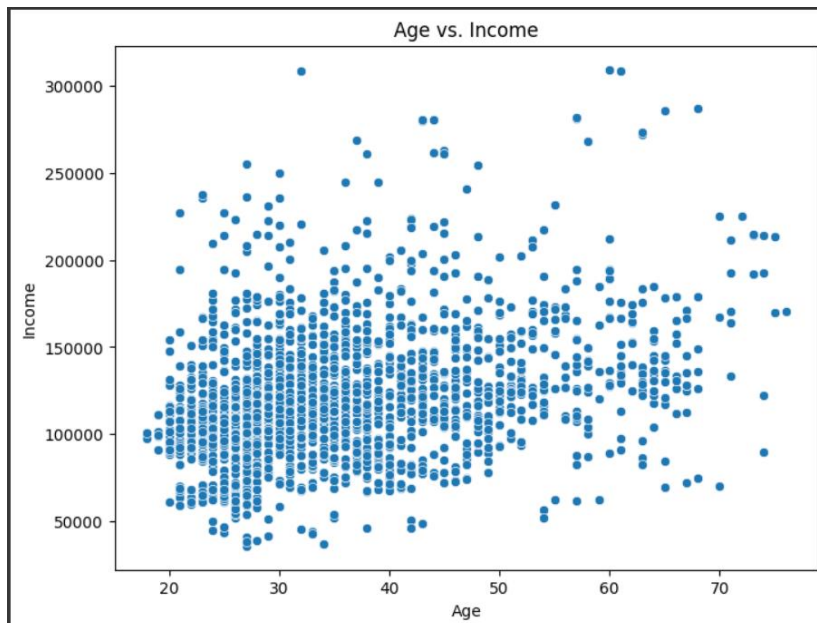
Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

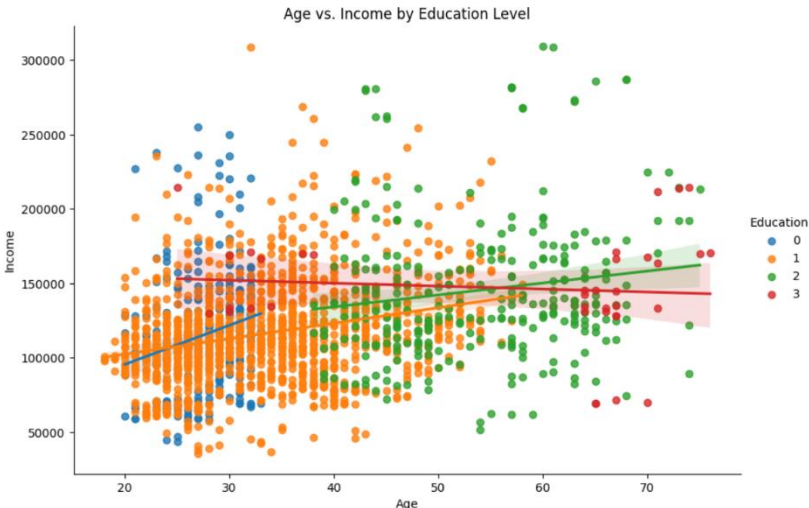
Section	Description																																																						
Data Overview	2000RowsX8Columns																																																						
	<table><tr><th></th><th>ID</th><th>Sex</th><th>Marital status</th><th>Age</th><th>Education</th><th>Income</th><th>Occupation</th><th>Settlement size</th></tr><tr><td>0</td><td>100000001</td><td>0</td><td>0</td><td>67</td><td>2</td><td>124670</td><td>1</td><td>2</td></tr><tr><td>1</td><td>100000002</td><td>1</td><td>1</td><td>22</td><td>1</td><td>150773</td><td>1</td><td>2</td></tr><tr><td>2</td><td>100000003</td><td>0</td><td>0</td><td>49</td><td>1</td><td>89210</td><td>0</td><td>0</td></tr><tr><td>3</td><td>100000004</td><td>0</td><td>0</td><td>45</td><td>1</td><td>171565</td><td>1</td><td>1</td></tr><tr><td>4</td><td>100000005</td><td>0</td><td>0</td><td>53</td><td>1</td><td>149031</td><td>1</td><td>1</td></tr></table>		ID	Sex	Marital status	Age	Education	Income	Occupation	Settlement size	0	100000001	0	0	67	2	124670	1	2	1	100000002	1	1	22	1	150773	1	2	2	100000003	0	0	49	1	89210	0	0	3	100000004	0	0	45	1	171565	1	1	4	100000005	0	0	53	1	149031	1	1
		ID	Sex	Marital status	Age	Education	Income	Occupation	Settlement size																																														
	0	100000001	0	0	67	2	124670	1	2																																														
	1	100000002	1	1	22	1	150773	1	2																																														
	2	100000003	0	0	49	1	89210	0	0																																														
	3	100000004	0	0	45	1	171565	1	1																																														
4	100000005	0	0	53	1	149031	1	1																																															

Univariate Analysis



Bivariate Analysis



Multivariate Analysis																																																							
Outliers and Anomalies	-																																																						
Data Preprocessing Code Screenshots																																																							
Loading Data	<table><tr><th></th><th>ID</th><th>Sex</th><th>Marital status</th><th>Age</th><th>Education</th><th>Income</th><th>Occupation</th><th>Settlement size</th></tr><tr><td>0</td><td>100000001</td><td>0</td><td>0</td><td>67</td><td>2</td><td>124670</td><td>1</td><td>2</td></tr><tr><td>1</td><td>100000002</td><td>1</td><td>1</td><td>22</td><td>1</td><td>150773</td><td>1</td><td>2</td></tr><tr><td>2</td><td>100000003</td><td>0</td><td>0</td><td>49</td><td>1</td><td>89210</td><td>0</td><td>0</td></tr><tr><td>3</td><td>100000004</td><td>0</td><td>0</td><td>45</td><td>1</td><td>171565</td><td>1</td><td>1</td></tr><tr><td>4</td><td>100000005</td><td>0</td><td>0</td><td>53</td><td>1</td><td>149031</td><td>1</td><td>1</td></tr></table>		ID	Sex	Marital status	Age	Education	Income	Occupation	Settlement size	0	100000001	0	0	67	2	124670	1	2	1	100000002	1	1	22	1	150773	1	2	2	100000003	0	0	49	1	89210	0	0	3	100000004	0	0	45	1	171565	1	1	4	100000005	0	0	53	1	149031	1	1
	ID	Sex	Marital status	Age	Education	Income	Occupation	Settlement size																																															
0	100000001	0	0	67	2	124670	1	2																																															
1	100000002	1	1	22	1	150773	1	2																																															
2	100000003	0	0	49	1	89210	0	0																																															
3	100000004	0	0	45	1	171565	1	1																																															
4	100000005	0	0	53	1	149031	1	1																																															
Handling Missing Data	<pre>df.isnull().sum()</pre>																																																						
Data Transformation	<pre>from sklearn.preprocessing import StandardScaler sc=StandardScaler() x=sc.fit_transform(x)</pre>																																																						
Feature Engineering	<pre># Feature Engineering - One-hot encode categorical variables (Sex, Marital status, Education, Occupation, Settlement size) df = pd.get_dummies(df, columns=['Sex', 'Marital status', 'Education', 'Occupation', 'Settlement size'], drop_first=True)</pre>																																																						

Save Processed Data

```
import pickle
pickle.dump(xgb_model, open("xgbmodel.pkl", 'wb'))

from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
pickle.dump(sc, open('scaling.pkl', 'wb'))
```