

Hi Sprocket Central Pty Ltd ,

myself , Yogesh Gavali from KPMG team. I am providing your esteemed company the details about the quality of the data which has been shared.

The data has been assessed carefully and the following mitigations are proposed. Please have a look and lets us know if there are any concerns .

Before proceeding further , the preliminary assessment of the datasets is as follows:

	Transactions Data	Customers data
Number of Observations	Transactions :20000	Customer Demographic : 4000 Customer Address : 3999
Number of columns	Transactions :13	Customer Demographic : 13 Customer Address : 6

Based on the preliminary assessment the Customer data can be merged into one by combining existing and new customer list as :

Customer Data:

- i. The New dataset ‘Customer’ is created by merging ‘CustomerDemographic’ and ‘CustomerAddress’.

Transaction Data :

- i. We can see that customer_id has ranging from 1-3500. An outlier of customer_id of ‘5034’.

To handle null values the missing rows are imputed with ‘Unknown’ or the rows are removed as they contribute less than 1% of total data.

		Customers Data	Transactions data
Missing values	Observed	<div><div>1. last_name125</div><div>2. DOB87</div><div>3. job_title506</div><div>4. job_industry656</div><div>5. tenure87</div><div>6. address4</div><div>7. postcode4</div><div>8. state4</div><div>9. country4</div><div>10. property_valuation4</div></div>	<div><div>1. online_order360 brand</div><div>2. 197 product_line197</div><div>3. product_class197</div><div>4. product_size197</div><div>5. standard_cost197</div><div>6. product_first_197</div><div>7. sold_date</div></div>
	Mitigation	<div><div>1. last_nameFilled with ‘Unknown’</div><div>2. DOBconverted to age, age_group</div></div>	<div><div>1. online_orderDropped null rows</div><div>2. brandFilled with ‘Unknown’</div><div>3. product_lineFilled with ‘Unknown’</div><div>4. product_classFilled with ‘Unknown’</div></div>

Quality Check :

		<p>and imputed with mode value.</p> <p>3. job_title Filled with 'Unknown'</p> <p>4. job_industry_category Filled with 'Unknown'</p> <p>5. tenure Filled with mean value</p> <p>6. address Filled with 'Unknown'</p> <p>7. postcode Dropped null rows</p> <p>8. state Dropped null rows</p> <p>9. country Filled with 'Unknown'</p> <p>10. property_valuation Dropped null rows</p>	<p>5. product_size Filled with 'Unknown'</p> <p>6. standard_cost Dropped null rows</p> <p>7. product_first_sold_date Dropped null rows</p>
Outlier values	Observed	1. DOB has outlier of '1843-1221'	1. customer_id has value of 5034
	Mitigation	1. This is row is removed.	1. This is row is removed.
Inconsistence values	Observed	<p>1. State – (NSW,VIC,QLD, New South Wales, Victoria)</p> <p>2. Gender- ('F', 'Femal', 'M', 'U')</p>	
	Mitigation	<p>1. Replaced 'M': 'Male', 'Femal' : 'Female', 'F' : "Female", 'U' : "Unknown"</p> <p>2. Replaced 'New South Wales': 'NSW', 'Victoria': 'VIC'</p>	
Inconsistence data type	Observed		1. product_first_sold_date is in float64 type
	Mitigation		1. product_first_sold_date is converted to datetime
Columns dropped	Observed	<p>1. Default</p> <p>2. deceased_indicator</p>	
	Mitigation	1. These columns are dropped as they are used in calculating 'Rank' column.	

		<p>2. 'default' Column contains metadata and does not provide any additional information</p> <p>3. deceased_indicator is dropped as we have 2 rows in Y category , It doesn't provide much insight.</p>	
<i>Duplicated Values</i>	Observed	1. No duplicate rows	1. No duplicate rows
	Mitigation		

Please provide the acknowledgement for the assumptions taken while this data quality assessment so we can proceed for the data cleaning and further analysis required for modelling.

Regards,

Yogesh Gavali

(KPMG team)