

## Lead Scoring Case Study: Summary Report

Process Followed	<ol style="list-style-type: none"> <li>1. Initial brainstorming &amp; approach formulation: <ul style="list-style-type: none"> <li>➤ Involved walkthrough of the original dataset, data dictionary and gaining some domain expertise followed by draft formulation of the approach to be followed.</li> </ul> </li> <li>2. upGrad session for validating the approach: <ul style="list-style-type: none"> <li>➤ The session helped in better understanding of the expectations from the case study and the business aspects.</li> </ul> </li> <li>3. Alignment on coding protocols and Data preparation <ul style="list-style-type: none"> <li>➤ Team documented the final steps &amp; prepared a common coding skeleton.</li> <li>➤ Data preparation involved reshaping the data based on business requirements. Some steps involved in that were: Dropping Columns/Rows, Handling Null Values, Outlier Treatment and EDA to gather insights from the data.</li> <li>➤ Precise coding along with relevant comments helped in the clarity of the steps.</li> </ul> </li> <li>4. Modeling and validations: <ul style="list-style-type: none"> <li>➤ First followed RFE approach to pick around 25 variables for the first model followed by dropping 1 variable at a time, based on high p-value/VIF – using statsmodels.</li> <li>➤ The final model had 14 variables each having p-value&lt;0.05 and VIF&lt;2.5(indicates almost negligible multicollinearity).</li> <li>➤ Various metrics to measure the efficiency of model was calculated and then the same were evaluated on test data.</li> <li>➤ Team reviewed each other's coding work and fine tuned</li> </ul> </li> <li>5. Final consolidation: <ul style="list-style-type: none"> <li>➤ Train: Accuracy:0.797, Sensitivity:0.847, Specificity:0.766, Precision:0.696</li> <li>➤ Test: Accuracy:0.801, Sensitivity: 0.854, Specificity: 0.769, Precision:0.694</li> <li>➤ Lead Score calculated for the complete data and visualizations done to check the relation between Converted and Lead Score (High Lead Score mostly corresponds to Converted value 1)</li> </ul> </li> </ol>
Learnings	<ol style="list-style-type: none"> <li>1. Handling categorical variables with numerous categories <ul style="list-style-type: none"> <li>➤ Since we do not want our dataset to have huge number of dummy variables, it better to group categories having less weightage into 1 single category as "Others".</li> </ul> </li> <li>2. Handling skewed categorical variables</li> </ol>

	<ul style="list-style-type: none"> <li>➤ Highly skewed categorical columns indicate its less impact to the overall outcome. Hence such columns should be dropped</li> </ul> <p>3. NULL imputation in categorical variables</p> <ul style="list-style-type: none"> <li>➤ Null imputation by mode would not always be the best solution for categorical columns.</li> </ul> <p>4. Process of finding the right cut-off</p> <ul style="list-style-type: none"> <li>➤ Checking the Predicted value against multiple cut-off values is a manual approach. Can be done by plotting accuracy, sensitivity and specificity for various probabilities. The intersection of the graphs gives the optimum cut-off. Anything less than that would increase False Positives leading to increase in Sensitivity but decrease in Specificity and anything above that would increase Specificity but decrease Sensitivity.</li> </ul> <p>5. Adjusting the predictions based on business needed</p> <ul style="list-style-type: none"> <li>➤ The scenario-based questions helped us think differently on how to make the business grow in a better way.</li> </ul>
--	---