# Lead Scoring Case Study Summary

**Problem Statement:**

X Education is an education company that offers online courses for industry professionals. The company attracts many visitors to its website through various marketing channels. The company faces a problem: its lead conversion rate is very low. Out of 100 leads, only 30 become customers on average. To solve this problem, X Education wants to identify the most potential leads, also known as 'Hot Leads'. The company has hired you to help them with this task. Your job is to build a model that can assign a lead score to each lead based on various factors, such as their demographics, behavior, preferences, etc. The higher the lead score, the more likely the lead is to convert. The lower the lead score, the less likely the lead is to convert. The company's CEO has set a target of achieving an 80% lead conversion rate with this model.

**Solution Summary:**

1. **Step 1: Reading and Understanding Data**

   o   The data was loaded and thoroughly reviewed for analysis.

2. **Step 2: Data Cleaning**

   o   Variables with unique values were removed.

   o   Columns with 'Select' values (indicating no selection) were converted to null.

   o   Columns with more than 52% missing data were dropped, except for "Lead Quality", which was kept and imputed with "Not Sure" due to its significance.

   o   The data was cleaned by eliminating skewed or duplicate variables, filling missing values for numerical features, and creating new categories for categorical variables.

   o   Inconsistent labels (upper vs. lowercase) were standardized.

   o   Sales team-generated features were removed to avoid confusion.

3. **Step 3: Data Transformation**

   o   Binary variables were converted to '0' and '1'.

4. **Step 4: Creating Dummy Variables**

   o   Dummy variables were created for categorical features, and redundant ones were removed.

5. **Step 5: Test-Train Split**

   o   The dataset was split into training and testing sets with a 70-30% distribution.

6. **Step 6: Feature Rescaling**

- o All variables were scaled using Standard Scaling.

- o A heatmap was plotted to check for correlations between variables.

7. **Step 7: Model Building**

- o **Feature Selection**: Recursive Feature Elimination (RFE) was used to select the 15 most important features.

- o Insignificant features were removed based on P-values, leaving 12 key features with good Variance Inflation Factors (VIF).

- o The optimal probability cutoff was determined by evaluating accuracy, sensitivity, and specificity.

- o An ROC curve was plotted, achieving 95% area coverage, confirming the model's effectiveness.

- o The model was evaluated to ensure it correctly predicted 80% of converted cases.

- o Precision, recall, accuracy, sensitivity, and specificity were calculated on the training set.

- o A cutoff value of 0.25 was chosen based on the precision-recall trade-off.

- o The model was tested on the test set, producing a conversion rate of 90.78%, with sensitivity at 84.12% and specificity at 94.58%.

8. **Step 8: Conclusion**

- o The model's conversion rate on the test set was 84%, surpassing the CEO's target of 80%.

- o The model's strong sensitivity helps identify the most promising leads.

- o The top features influencing lead conversion are:

  - ▪ Tags_Lost to EINS

  - ▪ Tags_Closed by Horizzon

  - ▪ Tags_Will revert after reading the email