# LEAD SCORING CASE STUDY:

**Approach**

1. Importing Data

2. Inspecting the Dataframe

3. Data Preparation (Encoding Categorical Variables, Handling Null Values)

4. EDA (univariate analysis, outlier detection, checking data imbalance)

5. Dummy Variable Creation

6. Test-Train Split

7. Feature Scaling

8. Looking at Correlations

9. Model Building (Feature Selection Using RFE, Improvising the model further inspecting adjusted R-squared, VIF and p-vales)

10. Build final model

11. Model evaluation with different metrics Sensitivity, Specificity

**Conclusion:**

- The logistic regression model predicts the probability of the target variable having a certain value, rather than predicting the value of the target variable directly. Then a cutoff of the probability is used to obtain the predicted value of the target variable.

- Here, the logistic regression model is used to predict the probability of conversion of a customer.

- Optimum cut off is chosen to be 0.27 i.e. any lead with greater than 0.27 probability of converting is predicted as Hot Lead (customer will convert) and any lead with 0.27 or less probability of converting is predicted as Cold Lead (customer will not convert)

- Our final Logistic Regression Model is built with 14 features.

- Features used in final model are ['Do Not Email', 'Lead Origin_Lead Add Form', 'Lead Source_Welingak Website', 'Last Activity_SMS Sent', 'Tags_Busy', 'Tags_Closed by Horizzon', 'Tags_Lost to EINS', 'Tags_Ringing', 'Tags_Will revert after reading the email', 'Tags_switched off', 'Lead Quality_Not Sure', 'Lead Quality_Worst', 'Last Notable Activity_Modified', 'Last Notable Activity_Olark Chat Conversation']

- The top three categorical/dummy variables in the final model are 'Tags_Lost to EINS', 'Tags_Closed by Horizzon', 'Lead Quality_Worst' with respect to the absolute value of their coefficient factors.

'Tags_Lost to EINS', 'Tags_Closed by Horizzon' are obtained by encoding original categorical variable 'Tags'. 'Lead Quality_Worst' is obtained by encoding the categorical variable 'Lead Quality'.

- Tags_Lost to EINS (Coefficient factor = 9.578632)

- Tags_Closed by Horizzon (Coefficient factor = 8.555901)

- Lead Quality_Worst (Coefficient factor =-3.943680)

- The final model has Sensitivity of 0.928, this means the model is able to predict 92% customers out of all the converted customers, (Positive conversion) correctly.

- The final model has Precision of 0.68, this means 68% of predicted hot leads are True Hot Leads.

- We have also built an reusable code block which will predict Convert value and Lead Score given training, test data and a cut-off. Different cutoffs can be used depending on the use-cases (for eg. when high sensitivity is required, when model have optimum precision score etc.)