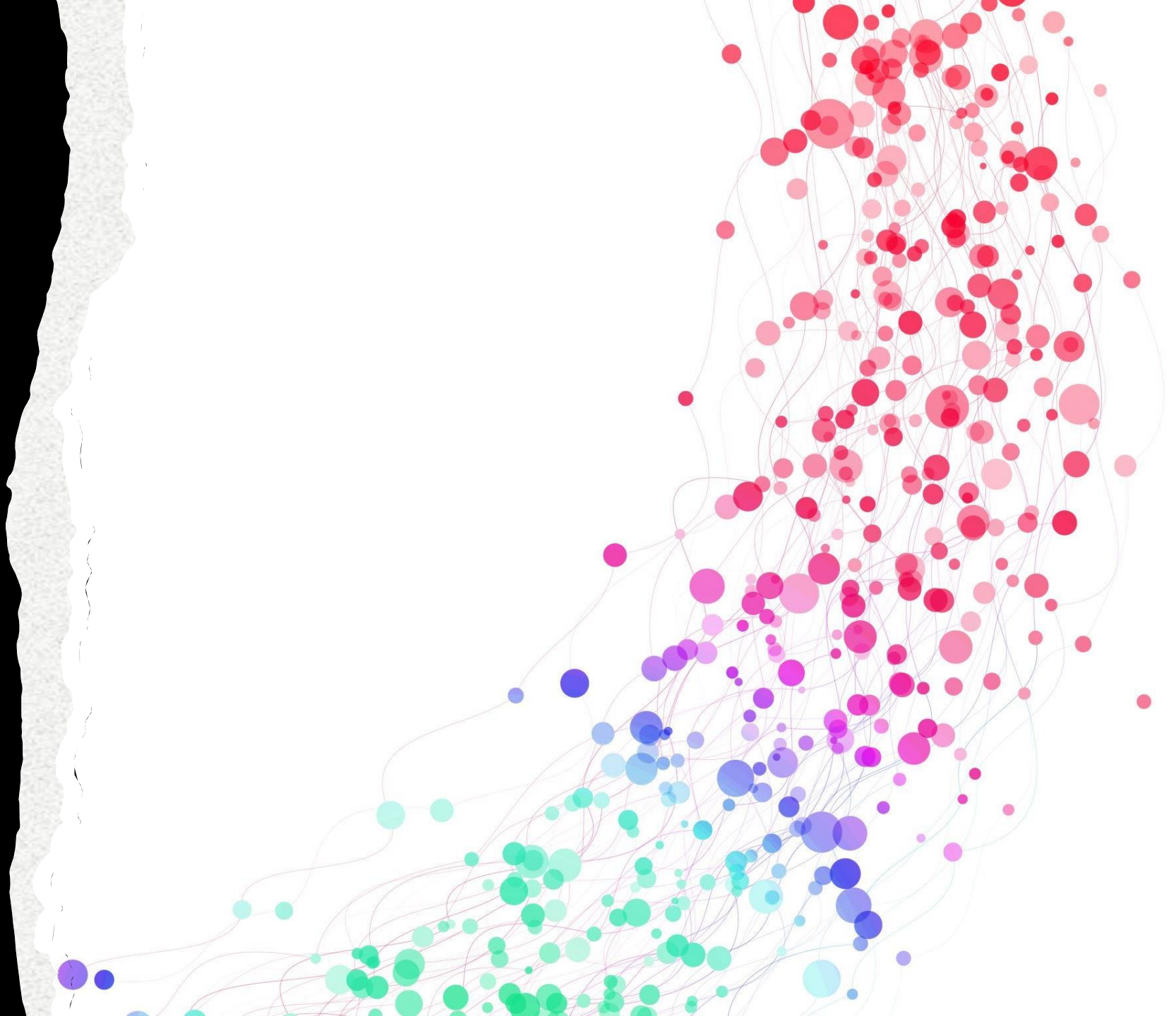# Lead Scoring

# Problem Description

- X Education is an education company that offers online courses to industry professionals

- The company advertises its courses on various websites and search engines, including Google

- X Education wants to improve its lead conversion rate by identifying the most promising leads - those who are most likely to become paying customers

Approach

Based on the problem description, it has been determined that the problem is a classification problem. As a result, logistic regression was chosen as the method for calculating the lead rate. The following steps were taken to solve the problem Data Reading and Understanding

- Number of rows and columns

- Data types of each columns

- Checking first few rows how data looks

- Checking how the data is spread

- Checking for duplicates, if any

# Data Cleaning

- Checking for any column names correction
- Checking for null values and imputing them with appropriate methods
  - We used mode imputation for categorical columns
  - We used mean imputation for numerical columns, if there is no skewness in data
  - We used median imputation for numerical columns, if there is skewness in the data
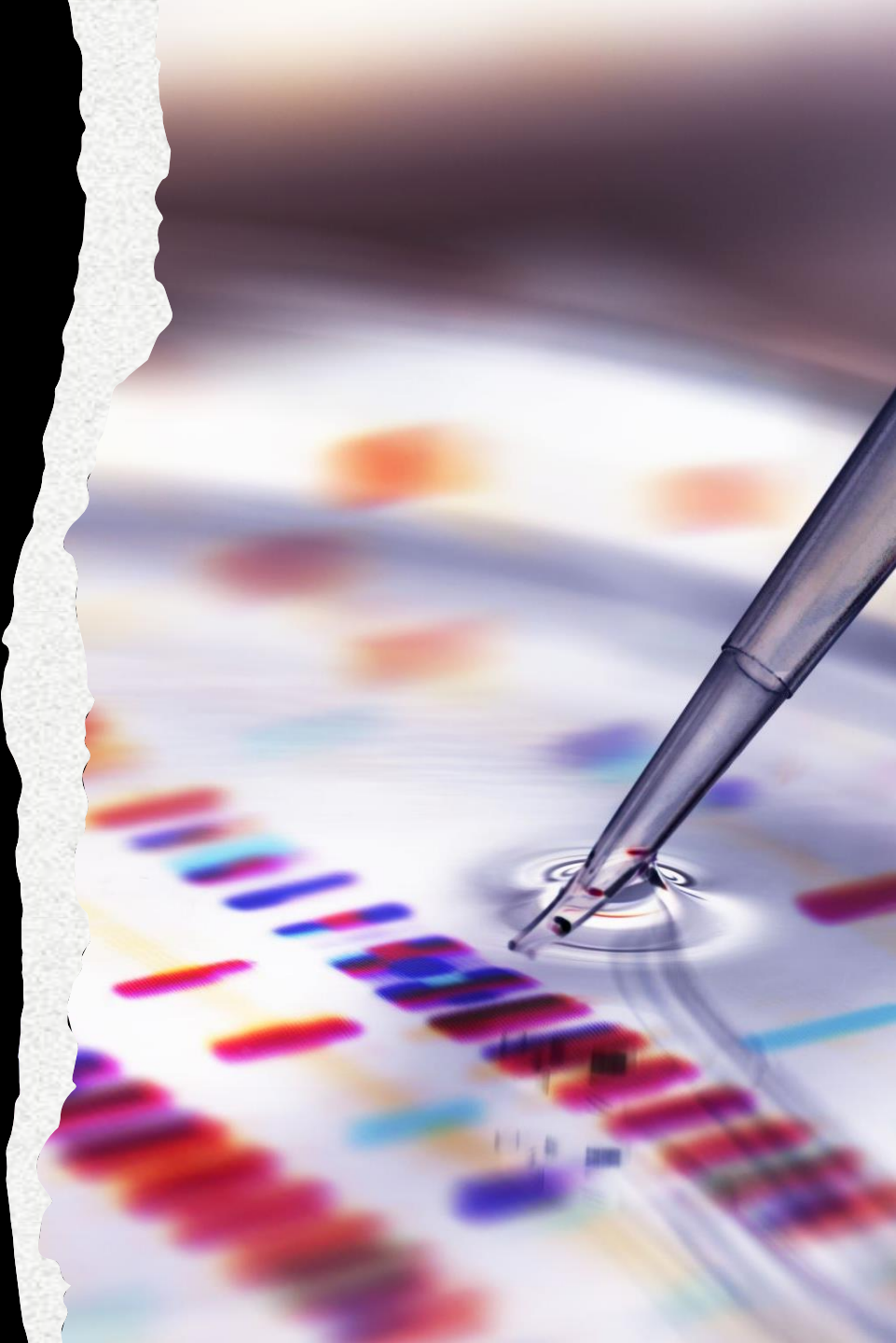
# Data Visualization and Outliers Treatment

- We conducted a univariate analysis on categorical columns to determine which columns were most relevant and removed those with near-zero variance

- We also performed a bivariate analysis on categorical columns to see how they varied with respect to the Converted column

- We used the IQR method to treat outliers in the dataset and plotted a correlation matrix to identify correlated columns
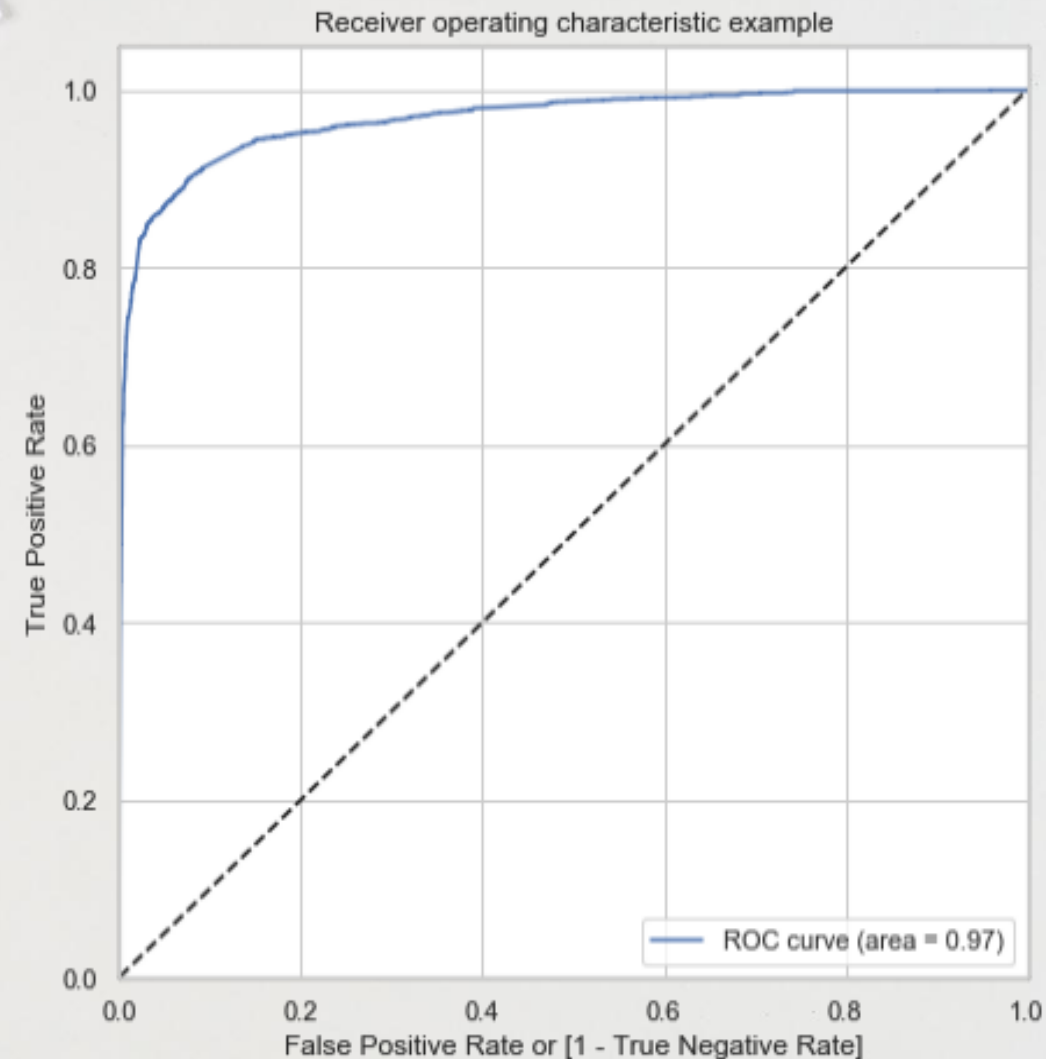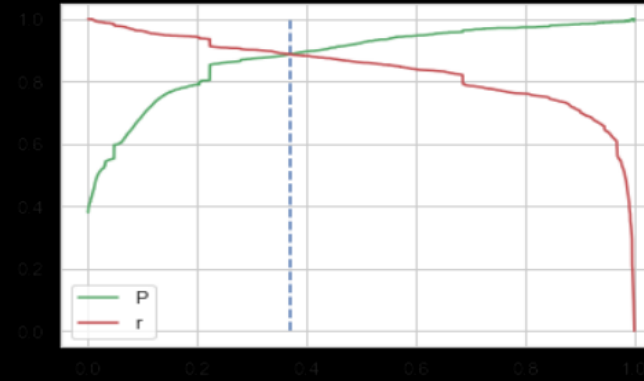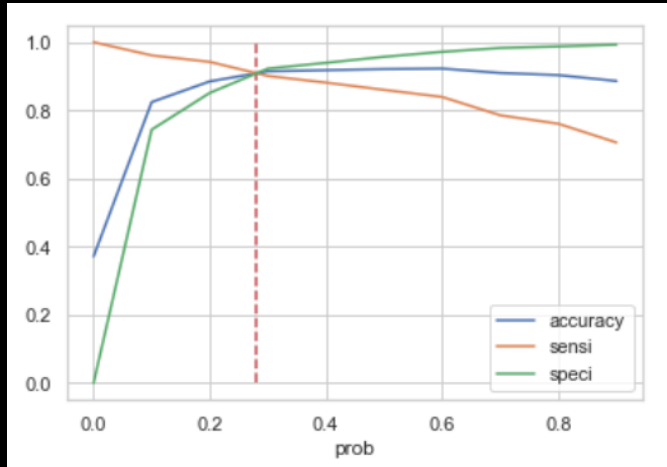
# Feature Scaling

At this point, our data was clean and free of outliers. Since logistic regression requires numerical input parameters, we converted all categorical columns to numerical values. Columns with only two levels, "Yes" and "No," were converted to numerical values using binary mapping. Columns with more than two levels were converted to dummy variables using the pd.get_dummies function. The resulting data contained only numerical columns and dummy variables. Before building the model, we rescaled all numerical columns using the standard scaler method

# Model Building

- We utilized the Recursive Feature Elimination technique to remove attributes and build a model on the remaining attributes

- RFE uses model accuracy to determine which attributes and combinations of attributes are most predictive of the target attribute

- After creating a stable model, we predicted probabilities on the train set and created a new "predicted" column with a value of 1 if the probability was greater than 0.5, and 0 otherwise



Receiver operating characteristic example. ROC curve (area = 0.97), True Positive Rate vs False Positive Rate or [1 - True Negative Rate].

# Model Evaluation on Train Set

- We calculated three metrics - accuracy, sensitivity, and specificity - for probabilities ranging from 0.0 to 0

- To make predictions on the train dataset, we determined the optimum cutoff to be 0.28 by finding the intersection of sensitivity, specificity, and accuracy

- For the test dataset, we considered the optimum cutoff to be obtained from the precision-recall graph of the train dataset

# Predictions on Test Set

After determining the optimum cutoff and calculating the metrics on the train set, we made predictions on the test dataset

Here are our observations: It seems that the model performed well on both the train and test datasets, with high accuracy, sensitivity, and specificity values

| Train Data | |
|---|---|
| Accuracy | 91.23% |
| Sensitivity | 90.53% |
| Specificity | 91.60% |

| Test Data | |
|---|---|
| Accuracy | 92.06% |
| Sensitivity | 90.53% |
| Specificity | 91.60% |

# Final Observations

- The model appears to accurately predict the conversion rate, which should enable the education company to select the most promising or "hot" leads

- There are certain variables that help drive hot leads