



# Sales & Mkt Summit: Optimizing Lead Conversion for X Education

Yuleica Ledezma  
Data Analyst, Corporate Sales and Marketing



# Executive Summary

- Objective:** Boost lead conversion rates by predicting high-conversion leads using machine learning models.
- Current Process:** Manual lead scoring based on intuition and limited data, resulting in inefficiency.
- Goal:** Develop a predictive model that improves the sales team's focus and efficiency, providing more precise leads, which will result after implementation (*outside the scope of this project*) to raising the current conversion rate from 38%.



# Performance Metrics

- Not only focusing on **Accuracy** (how often the model is correct)

**Precision:** Fewer false positives, enabling the sales team to focus on higher-quality leads.

- Recall:** Captured more potential leads, ensuring no opportunities were missed.

**Goal to reach 85% on each**



# Data Sources

From our partner department [Kaggle Leads Dataset](#), they created a dataset that comprises of:

- Website Activity
- Interest Survey/ Program Inquiry response
- Sales interactions
- Manual Scoring based on leads activity ( large set of missing values)

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	
1	Prospect I	Lead Num	Lead Orig	Lead Sour	Do Not Err	Do Not Ca	Converted	TotalVisi	Total Time	Page View	Last Acti	Country	Specializa	How did y	What is yc	What mat	Search	Magazine	Newspap	X Educat
2	7927b2df-	660737	API	Olark Cha	No	No	0	0	0	0	Page Visited on Web	Select	Select	Unemploy	Better Cari	No	No	No	No	No
3	2a272436-	660728	API	Organic Si	No	No	0	5	674	2.5	Email Ope	India	Select	Select	Unemploy	Better Cari	No	No	No	No
4	8cc8c611-	660727	Landing P	Direct Trai	No	No	1	2	1532	2	Email Ope	India	Business	Select	Student	Better Cari	No	No	No	No
5	0cc2d4f8-	660719	Landing P	Direct Trai	No	No	0	1	305	1	Unreacha	India	Media anc	Word Of M	Unemploy	Better Cari	No	No	No	No
6	32566b28-	660681	Landing P	Google	No	No	1	2	1428	U	V	W	X	Y	Z	AA	AB	AC	AD	
7	2058e6f0-	660680	API	Olark Cha	No	No	0	0	0	c	Newspap	Digital Ad	Through R	Receive M	Tags	Lead Qual	Update m	Get updat	Lead Profi	City
8	9fae7df4-	660673	Landing P	Google	No	No	1	2	1640	No	No	No	No	Interested	Low in Rel	No	No	Select	Select	02.Mediur
9	20ef72a2-	660664	API	Olark Cha	No	No	0	0	0	No	No	No	No	Ringin		No	No	Select	Select	02.Mediur
10	cfa0128c-	660624	Landing P	Direct Trai	No	No	0	2	71	No	No	No	No	Will revert	Might be	No	No	Potential	Mumbai	02.Mediur
11	a4f65dfc-	660616	API	Google	No	No	0	4	58	No	No	No	No	Ringin	Not Sure	No	No	Select	Mumbai	02.Mediur
12	2a369e35-	660608	Landing P	Organic Si	No	No	1	8	1351	No	No	No	No	Will revert	Might be	No	No	Select	Mumbai	02.Mediur
13	9bc8ce93-	660570	Landing P	Direct Trai	No	No	1	8	1343	No	No	No	No	Will revert	Might be	No	No	Select	Mumbai	01.High
14	8bf76a52-	660562	API	Organic Si	No	No	1	11	1538	No	No	No	No	Will revert	Low in Rel	No	No	Potential	Mumbai	02.Mediur
										No	No	No	No			No	No			02.Mediur
										No	No	No	No			No	No			02.Mediur
										No	No	No	No			No	No			02.Mediur
										No	No	No	No			No	No			02.Mediur
										No	No	No	No	Will revert	Might be	No	No	Select	Other Met	02.Mediur
										No	No	No	No	Lost to EINS		No	No	Select	Thane & C02	02.Mediur
										No	No	No	No	Will revert	Might be	No	No	Potential	Select	01.High
										No	No	No	No			No	No			02.Mediur

# Data Overview and Preprocessing

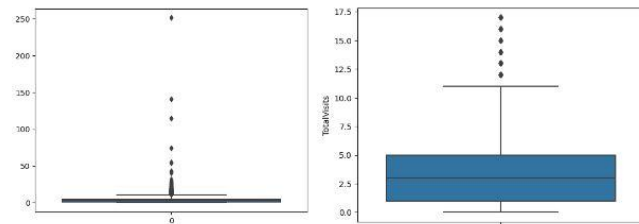
•Dataset Size: **9,241 rows and 37 columns**, sourced from a [Kaggle Leads Dataset](#).

- Numerical (7)
- Categorical (30)

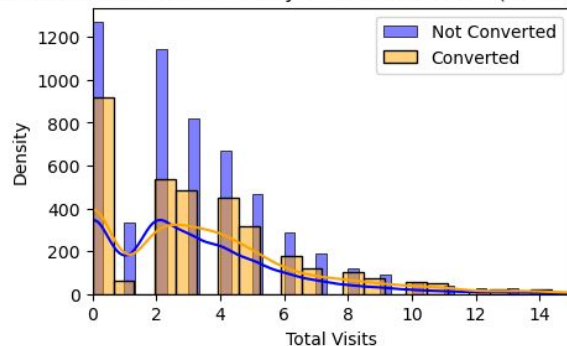
**Data Preprocessing:** Handle missing values, standardize categorical variables through one-hot encoding, and scale numerical features.

Challenges:

- Large number of null values (~15%, distributed unevenly in different columns)
- Skewness
- Outliers



Distribution of Total Visits by Conversion Status (Limited Range)

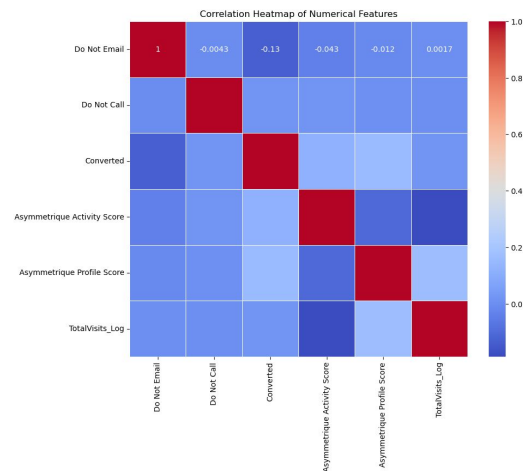
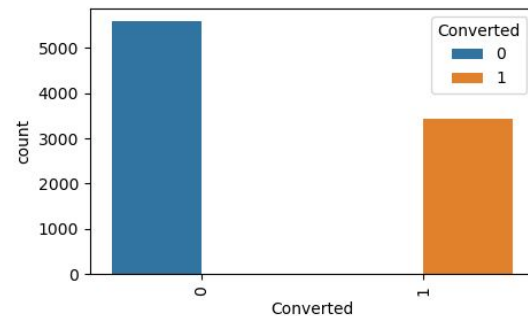


## Challenges (cont)

- Significant **class imbalance** (~ 60% not converted vs. ~40% converted), which could affect the models favoring the majority class. For this reason Accuracy is not the only metric, will focus also in Precision and Recall

- Affecting models
- Logistic Regression assumes a balanced class distribution

No strong correlations were found, allowing the model to use a broader set of features



# Machine Learning Models Tested

Models Evaluated:

- **Logistic Regression** (Baseline)

- Logistic Regression with PCA

- **Random Forest**

- Random Forest with Hyperparameter Tuning (GridSearchCV)

- **Neural Networks**

- Layers: 2 hidden layers

- Neurons:

- 64 neurons in the first layer

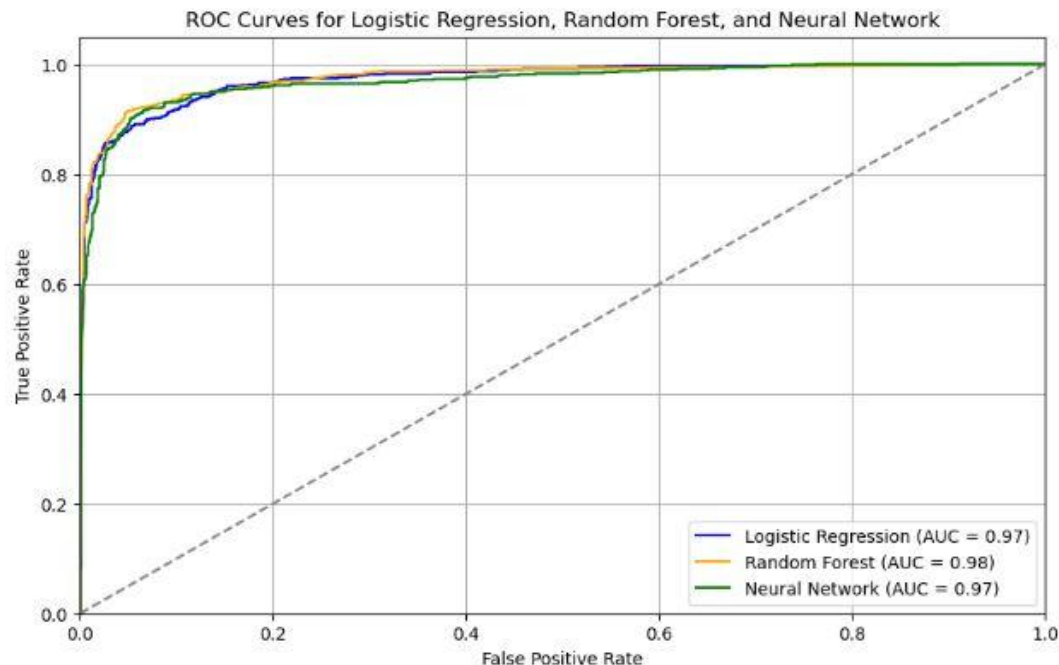
- 32 neurons in the second layer

- Activation Functions: ReLU for the hidden layers, Sigmoid for the output layer

- Neural Networks Keras Tuner

- **Performance Metrics:** Accuracy, Precision, Recall, and F1 Score.

# Model Evaluation and Results





# Model Evaluation and Results

	Log Regression(baseline)	Log Reg PCA	Random Forest	RF Hyperparameter Tuning	Deep Learning	Keras Tuner
Accuracy	.92	.92	.93	.94	.92	.94
Precision	.92	.93	.92	.94	.91	.95
Recall	.87	.87	.90	.90	.90	.89
F1 Score	.89	.90	.91	.92	.91	.91

# Recommendation

- Recommended Model: **Random Forest with Hyperparameter Tuning.**
  - Achieve a mix between **precision and recall**, ensuring the sales team focuses on the highest-potential leads without missing opportunities.

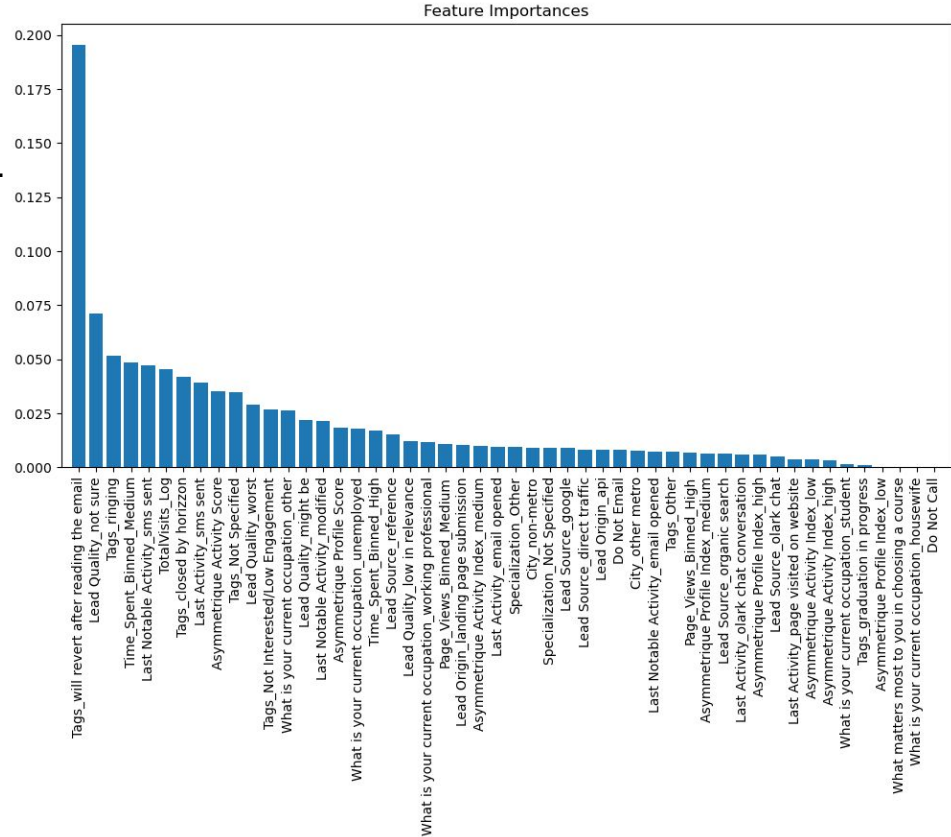
	Log Regression(baseline)	Log Reg PCA	Random Forest	RF Hyperparameter Tuning	Deep Learning	Keras Tuner
Accuracy	.92	.92	.93	.94	.92	.94
Precision	.92	.93	.92	.94	.91	.95
Recall	.87	.87	.90	.90	.90	.89
F1 Score	.89	.90	.91	.92	.91	.91

**Business Impact:** Saves time and boosts overall conversions.

# Random Forest:

- Works well with **large datasets**.
- Handles **non-linear relationships** effectively.
- **Robust to noise and outliers**.
- Provides **feature importance** insights.
- Handles **class imbalance** and **categorical/numerical data** well.
- **Hyperparameter tuning**, it offers flexibility to achieve high accuracy, precision, and recall.

Key features correlating with conversions:  
website engagement and last marketing  
activity.



# Conclusion

- Data-driven lead scoring improves efficiency and conversion rates.