**A REPORT**

**ON**

**Automating Lead Qualification Using ICP Matching and AI**

**BY**

| Names of the Students | ID No. |
|---|---|
| Divyam Gupta | 2023A7PS0423G |
| Ragav Krishna Ramesh | 2023A7PS0415G |

AT

CloudDefense.AI
A Practice School-I Station of

**BIRLA INSTITUTE OF TECHNOLOGY & SCIENCE, PILANI**

July, 2025

# A REPORT

## ON

# Automating Lead Qualification Using ICP Matching and AI

## BY

| Names of the Students | ID.No. | Discipline |
|---|---|---|
| Divyam Gupta | 2023A7PS0423G | Computer Science |
| Ragav Krishna Ramesh | 2023A7PS0415G | Computer Science |

Prepared in partial fulfillment of the
Practice School-I Course No.
BITS F221

## AT

CloudDefense.AI
A Practice School-I Station of



**BIRLA INSTITUTE OF TECHNOLOGY & SCIENCE, PILANI**

July, 2025

# Acknowledgment

We would like to express our sincere gratitude to **Dr. Sharan Gopal**, Faculty-in-Charge, for his constant support, encouragement, and valuable guidance throughout the course of our Practice School-I internship. His insights and mentorship played a crucial role in motivating us throughout the course of the project.

We are deeply thankful to our Project Mentor, **Mr. Varendra Maurya**, for his supervision and continuous assistance during the internship. His technical suggestions, practical insights, and problem-solving approach helped us tackle various challenges effectively.

We also wish to extend our gratitude to **Mr. Anshu Bansal**, for creating this invaluable opportunity for us and enabling a rewarding and enriching learning experience during our Practice School-I internship.

Their guidance was instrumental in the progress of our work, and we are truly grateful for the opportunity to learn and grow under their mentorship.

# BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE PILANI (RAJASTHAN)
## Practice School Division

**Station:** CloudDefense.AI          **Centre:** Online

**Date of Start:** 26/05/2025        **Date of Submission:** 16/07/2025

**Duration:** 2 Months

**Title of the Project**

Automating Lead Qualification Using ICP Matching and AI

**ID No./Names/Discipline of the students**

2023A7PS0415G – Ragav Krishna Ramesh – Computer Science

2023A7PS0423G – Divyam Gupta – Computer Science

**Name and designation of the expert**

Mr. Varendra Maurya – Head of Sales

**Name of the PS Faculty**

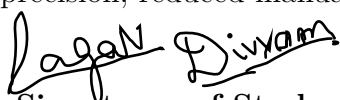Dr. Sharan Gopal

**Key Words**

Lead Qualification, Ideal Customer Profile, Automation, AI, Sales Funnel, Best Match Algorithm, Logistic Regression

**Project Areas**

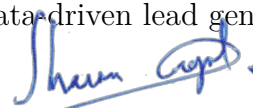Machine Learning, AI, Business Automation, Sales Tech

**Abstract**

This project focuses on automating lead qualification using AI techniques aligned with Ideal Customer Profiles (ICPs). The goal was to streamline sales outreach through machine learning models that analyze contact data and identify high-potential leads. Initially, we built a functional model combining AI and pattern recognition to match contacts with ICPs. We then refined it through Research and Development work by expanding ICP keywords and conducted a comparative analysis between logistic regression and neural networks to improve classification accuracy. The final system enhanced targeting precision, reduced manual effort, and enabled efficient, data-driven lead generation.

**Signatures of Students**          **Signature of PS Faculty**

Date: 16/07/2025          Date: 18-07-2025

# Table of Contents

# I. Introduction

This project aims to optimize and scale the outreach process for sellers onboarding onto our application across diverse digital platforms. In the modern data-driven landscape, pinpointing high-potential customers within vast datasets of professional and organizational information is both essential and challenging. To address this, we built a solution that utilizes a rich industry database to identify the most relevant leads, significantly enhancing targeting precision. Once qualified prospects are surfaced, the system leverages AI to craft and deliver personalized sales pitches through various communication channels, including LinkedIn, email, SMS, and cold calls. At its core, the solution features a machine learning pipeline centered on Ideal Customer Profiles (ICPs)—structured templates that define the most promising customers based on attributes like job role, industry, company size, location, revenue, and skills—enabling accurate and scalable lead qualification.

To operationalize this, we implemented a two-phase process: a Best Match algorithm first ranks ICPs based on similarity to predefined templates, ensuring strategic alignment with the product's ideal customer base. Following this, logistic regression models—trained on labeled datasets reflecting successful outreach—prioritize contacts most likely to convert. This integrated approach dramatically improves targeting precision and sales funnel efficiency. The pipeline is flexible and generalizable, allowing it to be adapted to varying product domains and sales objectives.

Since we had a functional model developed in the first phase, our primary objective in the second phase was to improve its accuracy through fine-tuning and research. We conducted an in-depth exploration of business and B2B-related terminology—such as engagement rate, funding stage, employee count, and tech stack—sourced from platforms like Kaggle, Crunchbase, and others. This research helped us redesign and optimize our Ideal Customer Profile (ICP) format to better reflect real-world lead qualification patterns and increase alignment with high-quality prospects.

To support large-scale testing and model training, we developed a CSV-based synthetic data generator that could simulate realistic contact datasets labeled according to our refined ICP logic. This tool allowed us to generate vast volumes of structured data efficiently. With this expanded dataset, we then performed a comparative analysis to determine the most effective AI model. We evaluated a neural network built using PyTorch against our earlier logistic regression-based model. The comparison involved assessing each model's ability to generalize across varied ICP patterns, handle imbalanced class distributions, and maintain consistency across different batches of synthetic and real-world data. This process enabled us to fine-tune our model architecture and data pipeline, laying the foundation for a more accurate and scalable lead qualification system.

# II. Work Plan and Methodology

## i. Research into Optimizing ICP Format

An **Ideal Customer Profile (ICP)** outlines the key characteristics of a company that is most likely to convert into a high-quality lead. These profiles are crucial for AI-based lead qualification, as they guide the model in identifying contacts that align closely with business goals. ICPs typically include dimensions such as industry vertical, company size, location, revenue range, technology usage, and decision-maker roles.

Our focus in this phase was to deepen the ICP structure by conducting extensive research into **B2B-specific keywords and attributes** that influence lead relevance. We analyzed real-world company data to identify important qualifiers such as *engagement rate, funding stage, employee count, IPO status, stock performance, technology stack, regulatory exposure*, and *pain points* commonly associated with enterprise tech adoption. This involved reviewing structured data across multiple sources to spot patterns in how successful B2B organizations are profiled and what terms frequently correlate with high-conversion leads.

After careful consideration and iterative refinement, we finalized the optimized ICP template format, and below is an example of an ICP in the final format:

```
{
    "industry": ["Healthcare Tech", "MedTech", "AI in Healthcare", "Wearable Tech"],
    "engagement_rate": "65-95",
    "company_size_employees": "100{800",
    "annual_revenue_usd": "10M{40M",
    "headquarters_location": "India",
    "technology_stack": ["Python", "AWS", "Kubernetes", "TensorFlow", "Edge AI"],
    "target_designations": ["Chief Medical Officer", "CTO", "Head of AI"],
    "pain_points": [
        "Medical device integration", "Data privacy compliance",
        "AI model explainability", "Real-time patient monitoring"
    ]
}
```

This refined template was essential in training the classification model, ensuring that it prioritized leads based on strong B2B alignment and sector-specific relevance.

# ii. Data Generation and Preprocessing

The dataset comprises various attributes related to professionals and their companies, such as job title, technologies used, industry, location, and financial metrics. The primary goal of preprocessing is to transform this raw data into a structured and learnable format for training machine learning models to classify Ideal Customer Profiles (ICPs).

## 1. Data Generation and Structuring

To support machine learning workflows for Ideal Customer Profile (ICP) classification, we first developed a data generation module capable of synthesizing structured datasets representing realistic professional and organizational profiles. This synthetic dataset included attributes such as *first name*, *last name*, *job title*, *company*, *seniority*, *department*, *company size (in terms of employees)*, *industry*, *relevant keywords*, *city*, *state*, *country*, *company address*, *company contact information*, *technology stack*, *annual revenue*, *total funding*, *latest funding amount*, *pain points*, and *engagement rate.*

Each entry was programmatically generated to align with specific ICP templates that defined the target audience across sectors such as Healthcare Tech and FinTech. These templates guided the assignment of contextual values to each field, simulating real-world diversity across regions, technologies, roles, and financial indicators. This dataset was critical for producing large volumes of labeled data suitable for training and testing classification models, especially in early phases when real labeled data was sparse or unavailable.

## 2. Data Preprocessing

After data generation, a structured preprocessing pipeline was applied to convert the raw, heterogeneous information into a high-quality feature space suitable for model training and evaluation.

### a. Feature Categorization
Profile data was organized into three key feature types for downstream processing:

- **Textual Features:** Cleaned and concatenated fields such as job title, department, industry, and location.

- **Numerical Features:** Quantitative attributes including employee count, revenue, funding amount, and engagement rate.

- **Set-Based Features:** Multi-label fields like technology stack, target roles, and pain points indicating categorical presence.

**b. Preprocessing Techniques**

To make the dataset suitable for machine learning models, we implemented a comprehensive preprocessing pipeline that handled both textual and numerical features. These transformations ensured consistency, reduced noise, and structured the raw data into a format amenable to learning algorithms.

- **Text Normalization:** All textual fields were lowercased, and special characters, punctuation, and redundant whitespace were removed. This step reduced vocabulary size and improved the reliability of subsequent vectorization techniques.

- **TF-IDF Vectorization:** Transformed cleaned textual data (e.g., job titles, industry, location) into numerical vectors by capturing the significance of words relative to the dataset. This emphasized informative terms while down-weighting frequent, less discriminative ones.

- **MultiLabel Binarization:** Applied to multi-valued categorical fields such as technologies and keyword tags. This technique represented the presence or absence of each possible category as binary features, allowing multiple attributes to be captured per profile.

- **One-Hot Encoding:** Used for single-valued categorical fields such as department or country. Each unique value was transformed into a distinct binary column, enabling categorical information to be included without imposing an ordinal structure.

- **Numerical Scaling:** Continuous attributes like annual revenue, employee count, total funding, and engagement rate were standardized using z-score normalization. This ensured uniform feature influence and improved convergence during training.

Together, these preprocessing techniques transformed the original heterogeneous dataset into a high-quality, learnable feature space. This enabled both classical and deep learning models to achieve higher accuracy and generalizability when identifying Ideal Customer Profiles.

## 3. Label Generation for ICP Compatibility

Training a model for direct lead-to-ICP matching ideally requires paired datasets, where each ICP is linked with a corresponding list of qualified leads. However, such labeled pairings are rarely available at scale and are difficult to generate manually. To overcome this limitation, we adopted a more scalable and flexible labeling strategy. Specifically, we transformed the dataset by introducing ten binary columns—each representing a predefined ICP template (e.g., ICP1, ICP2, ..., ICP10). A value of 1 in any column signifies that the lead matches the corresponding ICP, while a 0 indicates a mismatch. This reformulation enabled multi-label classification without relying on explicitly curated lead-ICP pairs and allowed us to leverage available data more effectively for model training.

# iii. Best Match Algorithm

The *Best Match algorithm* is designed to map an incoming lead profile to the most relevant Ideal Customer Profile (ICP) template. It does so by evaluating similarities across multiple attributes between the lead and each ICP. A composite score is calculated for each ICP, and the one with the highest similarity score is selected as the best fit.

## 1. Attribute-Wise Similarity Evaluation

To compare lead profiles against ICP templates, a set of key attributes is examined. Each attribute contributes a similarity score based on a suitable metric:

- **Industry Similarity:** Measured using the Jaccard similarity coefficient between the lead's and ICP's industry tags.

- **Company Size and Revenue Matching:** Evaluated using a *range overlap* score that quantifies the degree of intersection between value ranges.

- **Engagement Rate Compatibility:** Also computed using the range overlap score to assess how well the lead's engagement rate fits the ICP expectation.

- **Technology Stack Similarity:** Calculated using Jaccard similarity on the sets of technologies associated with both the lead and the ICP.

- **Designation Alignment:** Based on the overlap between target job titles from the ICP and the lead's actual designation.

- **Pain Point Match:** Jaccard similarity is used to measure the alignment between known pain points of the ICP and those indicated in the lead profile.

- **Geographic Location Match:** Evaluated with a binary score—*1 if the headquarters locations match, otherwise 0.*

## 2. Jaccard Similarity

For comparing sets like industry tags, tech stacks, and pain points, the Jaccard similarity is applied. Given two sets $A$ and $B$, it is defined as:

$$Jaccard(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

This value ranges from 0 (no overlap) to 1 (perfect match), providing an intuitive measure of set-based similarity.

### 3. Range Overlap Score

For numerical or bounded range features such as company size, revenue, and engagement rate, similarity is calculated using a normalized range overlap score:

$$Overlap(R_1, R_2) = \begin{cases} 0, & if \max(R_1^{low}, R_2^{low}) > \min(R_1^{high}, R_2^{high}) \\ \dfrac{\min(R_1^{high}, R_2^{high}) - \max(R_1^{low}, R_2^{low})}{\max(R_1^{high}, R_2^{high}) - \min(R_1^{low}, R_2^{low})}, & otherwise \end{cases}$$

A higher score indicates stronger alignment between the respective numerical ranges.

### 4. Final Scoring and Selection

Each feature contributes equally to the final compatibility score. The six computed scores—industry, size, revenue, engagement rate, tech stack, designation, pain points, and geography—are averaged to determine the best fit:

$$TotalScore = \frac{1}{8}\sum_{i=1}^{8} score_i$$

The ICP with the highest total score is returned as the most suitable match for the input profile.

## iv. Logistic Regression for ICP Classification

Logistic Regression plays a central role in our classification pipeline, serving as the baseline model to determine whether a given lead aligns with a predefined Ideal Customer Profile (ICP). As a probabilistic model designed for binary classification tasks, Logistic Regression outputs a probability value between 0 and 1, indicating the likelihood that a particular input corresponds to a positive class—in our case, whether a lead is a suitable match (label 1) or not (label 0) for the ICP.

To make this determination, the model learns a weighted combination of input features extracted during the preprocessing phase. These features, both textual (e.g., job title, industry) and numerical (e.g., company size, revenue, engagement rate), are transformed into a unified feature space using encoding techniques such as TF-IDF, binarization, and normalization. This ensures that all relevant attributes contribute proportionally to the final decision.

During training, the model receives input–output pairs: each input is a lead's feature

vector $\mathbf{x}$, and the output is a binary label $y \in \{0, 1\}$. The Logistic Regression model estimates the probability that a lead is a match using the sigmoid activation function:

$$P(y = 1 \mid \mathbf{x}) = \sigma(\mathbf{w}^T\mathbf{x} + b) = \frac{1}{1 + e^{-(\mathbf{w}^T\mathbf{x}+b)}}$$

Here, $\mathbf{w}$ is the vector of learned weights, $b$ is the bias term, and $\sigma(\cdot)$ is the sigmoid function that maps the linear combination of inputs to a probability.

The training objective is to minimize the binary cross-entropy loss, defined as:

$$\mathcal{L} = -\frac{1}{N}\sum_{i=1}^{N}[y_i\log(\hat{y}_i) + (1 - y_i)\log(1 - \hat{y}_i)]$$

where $\hat{y}_i$ is the predicted probability for the $i^{th}$ lead, and $y_i$ is the true label. This loss function penalizes incorrect predictions more severely as the predicted probability deviates from the true label, encouraging the model to be both confident and accurate. Overall, Logistic Regression provides a transparent and interpretable framework for ICP classification.

# v. Neural Network Architecture for ICP Classification

Neural networks offer greater flexibility in modeling non-linear relationships between input features and prediction targets. In our project, we designed a simple feedforward neural network named *ICPNet*, specifically tailored to predict lead compatibility with a given Ideal Customer Profile (ICP). Unlike linear models such as Logistic Regression, neural networks are capable of learning complex feature interactions and subtle patterns across diverse inputs.

The architecture of *ICPNet* consists of an input layer with 14 neurons, a hidden layer with 24 neurons activated by a ReLU function, and an output layer with a sigmoid activation for binary classification. The model is trained using the binary cross-entropy loss and optimized with the Adam optimizer.

```python
class ICPNet(nn.Module):
    def __init__(self):
        super(ICPNet, self).__init__()
        self.fc1 = nn.Linear(14, 24)
        self.relu = nn.ReLU()
        self.fc2 = nn.Linear(24, 1)
        self.sig = nn.Sigmoid()

    def forward(self, x):
        x = self.fc1(x)
        x = self.relu(x)
        x = self.fc2(x)
        x = self.sig(x)
        x = x.squeeze()
        return x
```

Before training, all numerical fields were scaled, and categorical fields were encoded using label encoding to convert the feature set into a numerical matrix. The dataset was then split into training and test sets, and tensors were created for input into the network.

During training, the model receives input tensors representing lead features and learns to minimize the binary cross-entropy loss between its predicted probability and the true label (match or non-match for a given ICP). Training is performed for 1000 epochs using backpropagation and Adam optimization.

Overall, this neural network complements the logistic regression model by capturing non-linear relationships that simpler models might overlook, making it a valuable addition to our classification pipeline.

# vi. Comparative Analysis: Logistic Regression vs Neural Network

In order to assess the effectiveness of different classification strategies for identifying Ideal Customer Profiles (ICPs), we implemented and evaluated two distinct approaches: Logistic Regression and a custom Neural Network model.

**1. Feature-Type Sensitivity**

Our dataset consists of a mix of numerical features (e.g., company size, revenue, funding) and set-based categorical fields (e.g., technologies used, keywords, departments). Logistic Regression showed slightly better performance on numerical attributes, especially when well-scaled and preprocessed. This is expected, as it is inherently a linear model that relies

heavily on the direct relationship between input features and the output probability. On the other hand, both models performed comparably on set-based fields (e.g., multi-label encoded technologies or pain points), where the importance of interactions is limited or sparsely distributed.

## 2. Model Flexibility and Generalization

A key difference lies in how the models generalize to unseen data. Logistic Regression is a static, interpretable model with coefficients representing linear weights for each feature. While this makes it highly explainable and fast to train, its capacity to capture complex, non-linear feature interactions is limited.

In contrast, Neural Networks offer a more dynamic and expressive representation. By passing input through multiple hidden layers and non-linear activations, neural models can learn richer patterns and interactions that are not linearly separable. This becomes particularly valuable when the dataset grows in size or when new ICP patterns emerge that were not strongly present in training data.

## 3. Training Behavior and Performance

Logistic Regression converges quickly and requires minimal tuning. It performs reliably on moderately-sized datasets with well-preprocessed features. Neural Networks, while slower to converge and more sensitive to hyperparameter tuning, show the potential to improve further when provided with more training epochs and data variety. In our experiment, both models achieved competitive accuracy, but logistic regression remained favorable due to its efficiency and robustness across folds.

## 4. Overall Insights

While neural networks are more expressive and adaptive, Logistic Regression proves to be a strong baseline, particularly for structured datasets with well-engineered features. It remains effective when interpretability, speed, and simplicity are priorities. However, in scenarios where customer behavior patterns are complex or evolve frequently, neural architectures offer better adaptability and performance headroom.

Hence, while both approaches have their merits, Logistic Regression remains a reliable model for initial ICP classification tasks, especially when combined with good preprocessing. Neural Networks, meanwhile, offer a more future-proof alternative for dynamic, non-linear environments where greater prediction depth is required.

# vii. Evaluation and Output Generation

Once trained, both the Logistic Regression and Neural Network (ICPNet) models are evaluated on a held-out test dataset to assess their classification performance. For each test sample, the model produces a probability score indicating the likelihood that a lead matches a given ICP template.

## 1. Prediction and Thresholding

During inference, the models output probability values in the range [0, 1]. These scores are converted into binary predictions using a threshold—typically set at 0.5. A prediction score $\geq 0.5$ is interpreted as a positive match (label 1), while a score $< 0.5$ is considered a non-match (label 0).

The performance of each model is quantitatively assessed using standard classification metrics. The primary metric used in our evaluation is:

$$Accuracy = \frac{Number\,of\,Correct\,Predictions}{Total\,Number\,of\,Predictions}$$

This metric provides a general estimate of how well each model is able to correctly classify leads across the test dataset.

## 2. Output Filtering and Lead Selection

Post-evaluation, both models generate a list of leads predicted to be positive matches for the specified ICP. These leads are filtered based on the predicted label and assembled into a usable output format.

For each selected lead, essential fields such as *first name*, *last name*, *title*, *company*, *city*, *state*, *country*, and *company phone* are extracted. This curated list serves as a high-confidence contact subset, which can then be exported or passed downstream for targeted outreach and business development.

By standardizing the evaluation and output generation pipeline across both models, we were able to conduct a fair comparison while ensuring the usability of results in a real-world lead qualification workflow.

# III. Results Achieved

## 1. Model Performance

Both Logistic Regression and Neural Network models were evaluated across multiple ICP templates to assess predictive accuracy. On average, the Logistic Regression classifier consistently achieved higher performance, with test accuracies ranging between **95% and 98%** across different ICPs. In comparison, the Neural Network model maintained a slightly lower but stable accuracy range of **85% to 90%**, reflecting its greater sensitivity to data volume and hyperparameter tuning. For example, in the case of *ICP4*, Logistic Regression reached a peak test accuracy of **98.00%**, while the Neural Network achieved **89.75%**, still within an acceptable decision threshold.

## 2. Best Match Effectiveness

The Best Match algorithm, responsible for assigning incoming leads to the most suitable ICP, demonstrated measurable improvement following ICP keyword expansion. Initially, for a representative input lead, the algorithm assigned it to *ICP4* with a similarity score of **0.86**, based on shared attributes across industry, company size, revenue range, and technology stack. However, after enriching the ICP definitions with additional contextual tags—including domain-specific keywords, emerging technologies, and updated professional titles—the same lead's similarity score rose to **0.92**. This uplift reflects a more precise semantic match and highlights the impact of continual keyword refinement on improving ICP associations. These enhancements increased both the confidence and interpretability of the lead-to-ICP mapping process.

## 3. Lead Retrieval Outcomes

Post-classification, the system returned high-quality leads matching both the predicted ICP and the input attributes. The retrieved contacts predominantly included decision-makers and senior professionals with high alignment across designation, industry focus, and company scale. This ensured a strong match for outbound targeting. The improved ICP keyword set and enriched input vectors contributed to a more precise lead retrieval process, thereby maximizing outreach relevance and minimizing cold-start mismatches.

# IV. Future Work

Moving forward, an important objective is to expand our collection of ICP templates to cover a broader range of industries, company sizes, and technical stacks. With a richer set of templates, we also need to generate and obtain high-quality labeled test data for each ICP to evaluate model performance more accurately and reliably. This will enable a more fine-grained understanding of how well the system generalizes across varying customer profiles.

Additionally, we aim to address the risk of overfitting by training our models on larger and more diverse datasets. Increasing the volume and variety of training data will improve robustness and generalizability.

The final step is integrating our classification pipeline into the broader project infrastructure, enabling seamless coordination and data flow with other teams working on enrichment, outreach, and analytics modules. This integration will help translate our model outputs into actionable business outcomes.

# V. Conclusion

The system successfully streamlines and automates the lead qualification pipeline by integrating machine learning, structured data engineering, and intelligent ICP modeling. It combines contextual lead matching, high-accuracy classification, and dynamic data generation to create a scalable outreach solution for sales and marketing teams.

- **Context-Aware Matching and Classification:** The Best Match algorithm maps leads to the most relevant ICPs using structured similarity. Expanding keyword sets improved match quality, while Logistic Regression offered strong performance on structured data. Neural Networks added flexibility for learning deeper patterns in complex feature spaces.

- **Targeted Contact Retrieval:** The pipeline surfaces high-potential leads filtered by predicted ICP match, role relevance, technological alignment, and geographic fit. This results in actionable, high-quality contact lists suitable for personalized outreach campaigns.

- **Synthetic Data Generation:** A custom data generator simulates realistic lead records across industries, roles, company sizes, and financial traits. This proved invaluable during the early development phase when real labeled data was scarce, and continues to support robust model training.

- **Scalable Preprocessing Pipeline:** A modular preprocessing architecture was built to convert raw inputs into high-quality features via TF-IDF vectorization, multi-label binarization, and numerical scaling. This consistent and flexible framework ensures compatibility with both classical and deep learning models.

- **Template-Driven ICP Modeling:** The ICP template architecture allows the system to adapt across industries by modifying key attributes such as revenue range, technologies, and pain points. This makes the platform easily extensible for multi-vertical use.

- **Automation and Efficiency Gains:** By automating the ICP matching and lead scoring processes, the system significantly reduces manual workload for business development teams. It enables them to concentrate on high-conversion targets, improving both productivity and pipeline precision.

- **Real-World Applicability:** The combination of keyword-rich ICP templates, accurate model predictions, and multi-channel delivery (LinkedIn, Email, SMS, etc.) positions the system for immediate adoption in real outbound sales environments.

In summary, this end-to-end lead qualification engine provides an intelligent, adaptive, and scalable solution for modern outreach workflows—balancing accuracy, interpretability, and extensibility in a production-ready framework.

# VI. Appendix

Below are the external resources and code notebooks used throughout the project:

- **Colab Notebook (Model Pipeline)**:
  https://colab.research.google.com/drive/1aHWkfwHBBKtoerOrOfls9nrhuYryT4hO?usp=sharingscr

- **Crunchbase API (Company Data)**:
  https://data.crunchbase.com/docs/getting-started

- **OpenCorporates (Legal Entity Information)**:
  https://opencorporates.com/

- **UK Government Company Info API**:
  https://developer.company-information.service.gov.uk/

- **Kaggle Datasets**:
  https://www.kaggle.com/

- **DBpedia (Structured Web Data)**:
  https://www.dbpedia.org/

- **Wikidata (Knowledge Base)**:
  https://www.wikidata.org/wiki/

- **PyTorch Documentation (Neural Network Implementation)**:
  https://pytorch.org/docs/stable/index.html

- **PyTorch Neural Network API Reference**:
  https://pytorch.org/docs/stable/nn.html

- **Google Colab Guide for PyTorch Users**:
  https://pytorch.org/tutorials/beginner/colab

- **Deep Learning Specialization by Andrew Ng (Coursera)**:
  https://www.coursera.org/specializations/deep-learning

- **CS231n: Convolutional Neural Networks for Visual Recognition (Stanford)**:
  https://cs231n.github.io/

- **Scikit-learn Documentation (for Logistic Regression)**:
  https://scikit-learn.org/stable/modules/linear$_m$odel.htmllogistic − regression

- **Evaluation Metrics for Classification Models**:
  https://scikit-learn.org/stable/modules/model$_e$valuation.html

# VII. References

[1] Suphakit Niwattanakul, Jatsada Singthongchai, Ekkachai Naenudorn, and Supacha-nun Wanapu,
*Using of Jaccard Coefficient for Keywords Similarity*,
Proceedings of the International MultiConference of Engineers and Computer Scientists 2013 (IMECS 2013), Vol. I, pp. 380–384.

[2] Kirill Eremenko and Hadelin de Ponteves,
*Machine Learning A-Z™: Hands-On Python & R In Data Science*,
Udemy Online Course.
https://www.udemy.com/course/machinelearning/

[3] Ilona Pawełoszek and Jerzy Korczak,
*An Approach to Discovery of Customer Profiles*,
Lecture Notes in Business Information Processing, Vol. 268, pp. 88–99, December 2016.
Presented at: International Conference on Research and Practical Issues of Enterprise Information Systems.

[4] Nikhil Ketkar,
*Introduction to PyTorch*,
In: Deep Learning with Python, pp. 195–208, October 2017.
DOI: 10.1007/978-1-4842-2766-4$_1$2

# VIII. Glossary

- **Ideal Customer Profile (ICP):** A representation of the company's most valuable customer segments, defined by attributes like industry, size, revenue, and job roles. ICPs guide targeting by identifying the highest-conversion potential leads.

- **Lead Qualification:** The process of evaluating potential customers to determine their likelihood of becoming paying clients. This ensures outreach is focused on high-value prospects.

- **Machine Learning Pipeline:** A series of steps involving data preprocessing, modeling, and evaluation to train and deploy ML models. It enables automated and reproducible learning-based workflows.

- **Best Match Algorithm:** A custom matching algorithm that assigns leads to the most relevant ICPs based on feature similarity. It compares attributes like industry, revenue, and technology stack.

- **Logistic Regression:** A supervised learning model used for binary classification tasks. It predicts the probability of a lead being a match to a given ICP.

- **TF-IDF Vectorization:** Term Frequency-Inverse Document Frequency is a technique that converts textual data into numerical features. It captures the importance of words relative to all documents.

- **Jaccard Similarity:** A statistical measure that quantifies similarity between two sets. It is calculated as the size of the intersection divided by the size of the union.

- **Range Overlap Score:** A similarity metric for numeric ranges such as revenue or employee count. It quantifies the degree of intersection between two value intervals.

- **Binary Classification:** A type of machine learning task where the output is one of two possible classes (e.g., match or no match). It's used here to determine if a lead fits an ICP.

- **Cross-Entropy Loss:** A loss function used in classification problems to measure the error between predicted and actual classes. It penalizes confident but incorrect predictions more heavily.

- **Multilabel Binarization:** A preprocessing technique that represents multiple binary labels as columns of 0s and 1s. Useful for handling multi-category tag-like data such as tech stacks.

- **Feedforward Neural Network:** A type of artificial neural network where information flows in one direction—from input to output—without cycles.

- **ReLU (Rectified Linear Unit):** A non-linear activation function commonly used in neural networks. It outputs the input directly if it's positive, otherwise it outputs zero.

- **Sigmoid Activation:** An activation function that maps input values to a probability range between 0 and 1, suitable for binary classification tasks.

- **PyTorch:** An open-source machine learning framework used to implement and train our custom neural network model.

- **Data Generator:** A custom script that synthetically produces labeled examples of leads for training and testing at scale, helping overcome data availability limitations.

- **Enrichment:** The process of enhancing raw lead data with additional contextual information such as job role, company metrics, or engagement signals.

- **Technology Stack:** A collection of technologies, frameworks, and tools used by a company. It plays a key role in evaluating ICP fit in tech-based B2B outreach.

- **Engagement Rate:** A performance metric indicating how actively a company or user interacts with products or content. Used to qualify lead potential.

- **Funding Stage:** A descriptor of a company's financial lifecycle (e.g., Seed, Series A, Series B). It provides insight into growth and investment status for lead evaluation.

- **IPO (Initial Public Offering):** The process by which a private company becomes publicly traded. Often used as a financial qualifier in ICPs.

- **Stock Performance:** A metric indicating how a publicly listed company's share value has changed over time. Used in evaluating financial maturity in B2B settings.