

**Email Classification for Support Team**

## **Email Classification for Support Team**

Prepared by: Siddharthan R

Email: [siddharthanr25@gmail.com](mailto:siddharthanr25@gmail.com)

Date: 23rd April 2025

# Email Classification for Support Team

## 1. Introduction

In today's digital age, organizations receive thousands of customer emails that need to be sorted, categorized, and processed efficiently. These emails may contain sensitive personal information, requiring secure handling to ensure user privacy.

This project focuses on automating the classification of emails into meaningful categories and detecting Personally Identifiable Information (PII) using state-of-the-art Natural Language Processing (NLP) techniques. It reduces manual effort, ensures data privacy, and enhances the overall workflow of support teams.

## 2. Objective

- Automatically classify emails into categories like Incident, Request, Problem, or Change.
- Detect and mask PII such as names, email addresses, phone numbers, and other sensitive entities.
- Provide a simple, fast, and secure API and UI interface to access these functionalities.

## 3. Methodology

The project uses a pipeline consisting of PII detection, masking, embedding generation using SBERT, and classification.

PII Detection:

A hybrid approach was used. spaCy's NER detected names, while regular expressions identified emails, phone numbers, card numbers, and other entities. Overlapping results were prioritized by match length.

PII Masking:

Sensitive data was replaced with placeholders like [email], [full\_name], and [phone\_number] to preserve the

## Email Classification for Support Team

message's context while ensuring privacy.

Embedding Generation:

Used all-mpnet-base-v2 from Sentence-BERT to convert masked text into semantic vectors, which outperform traditional TF-IDF embeddings.

Classification Models:

Trained three classifiers on SBERT embeddings: Logistic Regression, Random Forest, and XGBoost.

Evaluated using accuracy and F1-score.

### 4. Results

Dataset: 24,000 emails across four classes

Logistic Regression: Accuracy 69.5%, F1 Score 0.70

Random Forest: Accuracy 73.2%, F1 Score 0.68

XGBoost: Accuracy 77.0%, F1 Score 0.76 (Selected for deployment)

### 5. Challenges & Solutions

Overlapping Entities:

Solved by prioritizing longer matches during masking.

International Phone Numbers:

Enhanced regex patterns to support various global formats.

## Email Classification for Support Team

### Full Name Detection:

Upgraded spaCy model to en\_core\_web\_md for improved accuracy.

### Training Time:

Switched from traditional BERT models to SBERT + classical ML models for faster training without loss in performance.

### Performance Bottlenecks:

Cached SBERT embeddings to reduce latency in repeated operations.