# Email Classification and PII masking – Project Report

## Project Overview:

Support teams handle a large volume of emails on a daily basis, often containing sensitive information like full names, email addresses, credit card numbers, Aadhaar numbers, and phone numbers. Manually reviewing and categorizing these emails is time-consuming and error-prone. This project proposes an automated system to classify emails into four categories: **"Change, Incident, Problem, and Request"** while ensuring PII is detected and masked effectively.

## Deployment Details:

Hugging Face Deployment URL :
https://huggingface.co/spaces/Midhran/Email-classification

Git URL:
https://github.com/midhran/Email_Classification

## Technical Approach:

➢ The solution was implemented as a RESTful API using FastAPI and deployed on Hugging Face Spaces. It consists of two key components:
➢ PII Masking: Utilizes spaCy's Named Entity Recognition (NER) and regex patterns to identify and replace sensitive information with placeholder tokens such as [email], [full_name], etc.
➢ Email Classification: Converts input text into TF-IDF vectors and applies machine learning models to predict the appropriate support category.

## Model Development and Evaluation:

The email classification component was trained using traditional machine learning models suited for high-dimensional sparse text data. The models tested include Logistic Regression, Random Forest, and Naive Bayes. Each model was evaluated on class-wise F1-scores and overall accuracy.

### Evaluation Metrics Summary:

➢ Logistic Regression (Final Model)
Accuracy: 0.76
F1 Scores: Change (0.86), Incident (0.77), Problem (0.41), Request (0.91)

➢ Random Forest
Accuracy: 0.74
F1 Scores: Change (0.75), Incident (0.78), Problem (0.16), Request (0.89)

➢ Naive Bayes
Accuracy: 0.73
F1 Scores: Change (0.80), Incident (0.73), Problem (0.46), Request (0.90)

**Model Justification:**

Logistic Regression was chosen as the final model due to its robustness, speed, and interpretability. It offered the most balanced performance across all categories, particularly excelling in the 'Problem' category, where other models significantly underperformed. Its suitability for sparse feature spaces like TF-IDF makes it an ideal choice for this classification task.

## Challenges and Solutions:

➢ Several challenges were encountered during development and deployment:
 Package Compatibility: Initial issues arose with mismatched versions of spaCy and pydantic. This was resolved by locking compatible versions (spaCy 3.7.2, pydantic 1.10.13).
➢ Deployment Errors: Hugging Face required specific entry points and strict FastAPI structure. The application was restructured accordingly.
➢ PII Detection Precision: Over-masking occurred initially. The masking logic was refined by tuning regex expressions and filtering out false positives from NER.

## API Endpoint and Output:

**Input JSON:**
```
{
  "email_body": "Hi, my name is Jane Doe and my email is jane.doe@example.com. I need help with billing."
}
```

**Output JSON:**
```
{
  "input_email_body": "...",
  "list_of_masked_entities": [
    {"position": [18, 26], "classification": "full_name", "entity": "Jane Doe"},
    {"position": [47, 67], "classification": "email", "entity": "jane.doe@example.com"}
  ],
  "masked_email": "Hi, my name is [full_name] and my email is [email]. I need help with billing.",
  "category_of_the_email": "Incident"
}
```

## Deliverables Summary:

- ➤ Code Implementation:
  - Python scripts: `pii_masker.py`, `models.py`, `api.py`.
  - All modules include docstrings and follow PEP8 standards.
  - Requirements file and trained model files included.

- ➤ Report:
  - Overview of problem and technical approach.
  - Justification for chosen model.
  - Evaluation and deployment steps

- ➤ Final Output:
  - Fully functioning API hosted on Hugging Face Spaces.
  - Verified format compliance with strict assignment requirements