# Akaike Assignment - Email Classification System Project Report

# Submitted by: Mohammad Khaja Moinuddin

# Date: 24-04-2025

---

## Objective

To develop a production-ready Email Classification API capable of:

- Detecting and masking Personally Identifiable Information (PII) in customer support emails

- Classifying emails into predefined support categories

- Delivering results via a fast, scalable API

---

## Dataset

- **Source**: `combined_emails_with_natural_pii.csv`

- **Total Records**: 24,000 emails

- **Columns**:

  - `email`: email body text

  - `type`: support category label

---

## Approach

### Text Preprocessing

- TF-IDF Vectorizer used for text-to-feature vector conversion (5000 max features, stopwords removed)

### Model Training

- RandomForestClassifier with 100 estimators

- Dataset split into 80% training / 20% testing

- Model saved as `email_classifier.pkl` using `pickle`

---

## PII Masking Logic

Regex-driven rule-based masking for:

- Email addresses

- Phone numbers

- Aadhar numbers

- Credit/Debit card numbers

- CVV codes

- Expiry dates

- Full names

- Dates of birth

**Masked PII replaced with** `[entity_type]` placeholders.
Masked entities logged with start-end position, classification type, and original value.

---

## API Implementation

- **Framework**: FastAPI

- **Main Endpoint**: `POST /classify_email`

- **API Documentation**: `/docs` (Swagger UI)

**API Request Example**

```
{

  "input_email_body": "Hello John Doe, your aadhar number is
1234 5678 9012."

}
```

**API Response Example**

```
{

  "input_email_body": "Hello John Doe, your aadhar number is 1234
5678 9012.",

  "list_of_masked_entities": [

    {"position": [6, 14], "classification": "full_name", "entity":
"John Doe"},

    {"position": [35, 49], "classification": "aadhar_num", "entity":
"1234 5678 9012"}

  ],

  "masked_email": "Hello [full_name], your aadhar number is
[aadhar_num].",

  "category_of_the_email": "Billing"

}
```

---

## Deployment

- **Platform**: Hugging Face Spaces (Docker SDK)

- **Deployed via Dockerfile running Uvicorn server on port 7860**

  **Live API URL**
  https://huggingface.co/spaces/moin0317/email_classifier

---

## Deliverables Summary

- Deployed API on Hugging Face Spaces ✅

- PEP8-compliant modular code ✅

- Regex-based masking ✅

- Production-trained model ✅

- API with strict JSON response format ✅

- Final report (this document) ✅

---

## Challenges Faced

- Managing regex accuracy for multiple PII types without LLMs

- Handling Hugging Face Docker SDK deployment quirks

- Managing large model file size limits on GitHub

---

## Future Improvements

- Integrate transformer-based classification (if permitted)

- Extend regex coverage for multilingual PII patterns

- Add role-based API authentication

- Implement rate limiting, API key access

---

## Conclusion

The API is fully deployed, passes evaluation criteria, and meets Akaike submission requirements. It reliably masks PII, classifies emails, and serves predictions via a live, scalable API endpoint.