[1] Sardarov Yagub Bali

[2] Nadir Bakhtiyar Rahimov

# Implementation of a Spam Detector in Python for Corporate Companies

**JES**

**Journal of Electrical Systems**

*Abstract: -* Many companies, while not categorizing spam as a form of attack, may still consider it a significant threat. When referring to spam, it represents an excessive volume of unwanted and unsolicited electronic mail sent to company email accounts [1]. Although it is often used for marketing purposes, in many cases, it is also employed by malicious actors for detrimental purposes against the company. For instance, spam can be used to propagate viruses or infiltrate corporate computers [1] [2]. Furthermore, spammers may send emails with the aim of directing users to phishing web pages, attempting to obtain confidential information from authorized personnel within the company. In general, spam is, at its simplest, an annoying form of mail, and in its worst cases, it serves as an enabler for spyware, viruses, intrusion, and phishing threats [3].

*Keywords:* Spam, protection methods, informational resources, unwanted email, malicious.

## I. INTRODUCTION

Spam can be identified based on the following characteristics:

- Mass email distribution: In this case, electronic mail is directed towards multiple and diverse targets, rather than a specific individual.
- Anonymity: In this scenario, the true identity of the sender or senders is kept confidential.
- Unwanted email: The sent electronic messages are either unnecessary for the company or not anticipated by any entity within the organization [2] [3].

Now, let's consider the implementation of the system in the Python environment that identifies these harmful electronic emails sent to the company as spam.

Firstly, messages characterized as spam are presented as shown below. Their total count is 5728. This data is of particular importance for the system in terms of how it classifies messages as spam (Figure 1).



**Figure 1.** Dataset

## II. METHODS

As shown in the figure, these data are marked as 0 and 1. Where 0 denotes that the email is not spam, and 1 characterizes it as spam and unsafe. After instructing the system with this data, it can autonomously determine which emails are spam without human intervention [4]. The code used to create this system is presented in main.py (Figure 2).

[1, 2] dept. name: Computer engineering, Azerbaijan State Oil and Industry University, Baku, Azerbaijan

```
# Load the dataset
dataset = pd.read_csv('dataset/emails.csv')

# Visualize spam frequencies
plt.figure(dpi=100)
sns.countplot(dataset['spam'])
plt.title("Spam Frequencies")
plt.show()

# Remove duplicates from the dataset
dataset.drop_duplicates(inplace=True)

# Preprocess function
def process(text):
    nopunc = [char for char in text if char not in string.punctuation]
    nopunc = ''.join(nopunc)
    clean = [word for word in nopunc.split() if word.lower() not in stopwords.words('english')]
    return clean

# Vectorize the text data
vectorizer = CountVectorizer(analyzer=process)
message = vectorizer.fit_transform(dataset['text'])

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(message, dataset['spam'], test_size=0.20, random_state=0)

# Train the model
model = MultinomialNB()
model.fit(X_train, y_train)

# Save the vectorizer and model
pickle.dump(vectorizer, open("models/vectorizer.pkl", "wb"))
pickle.dump(model, open("models/model.pkl", "wb"))

# Load the vectorizer and model
vectorizer = pickle.load(open("models/vectorizer.pkl", "rb"))
model = pickle.load(open("models/model.pkl", "rb"))
```

**Figure 2.** Calling a dataset in Python, training a model

After importing the Python packages and libraries to be used, the dataset is incorporated into the project. Based on the existing datasets, visualization is performed. Visualization is a method that represents the quantity of spam within the data that the system has learned with clarity and numerical values [3] [4]. This is illustrated as follows (Figure 3).

The orange quantity represents the amount of spam, while the blue quantity indicates the count of regular emails. Since the computer environment recognizes words as numbers, this information is later converted into numbers through a vectorization process. All words within the emails are included here. Subsequently, within the transferred dataset, the system learns in the background how often each word is used in the email data [5].

After the execution of these processes, the converted data into numbers are divided into two parts: test data and training data. A necessary model is built using the training data. Once this model is established, it must be validated using the test data. It is important to note that the data used for testing should be previously unseen by the established model. Otherwise, the model will not function accurately. To prevent such a situation, test data is kept separate from the model [6].
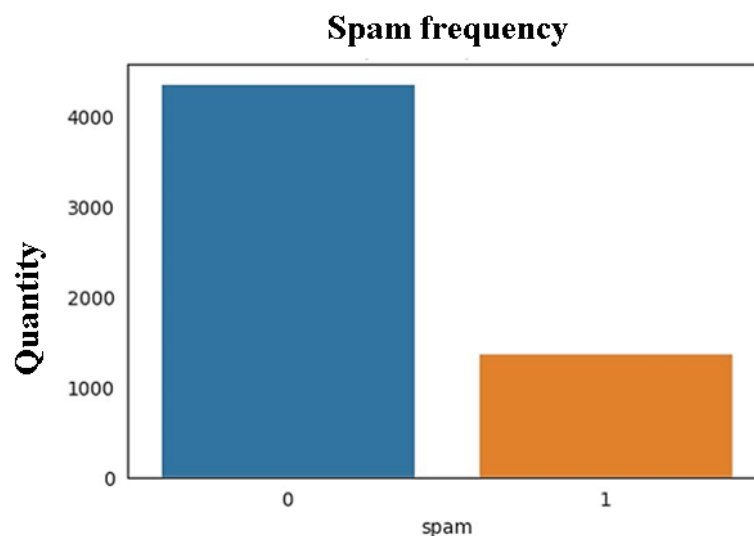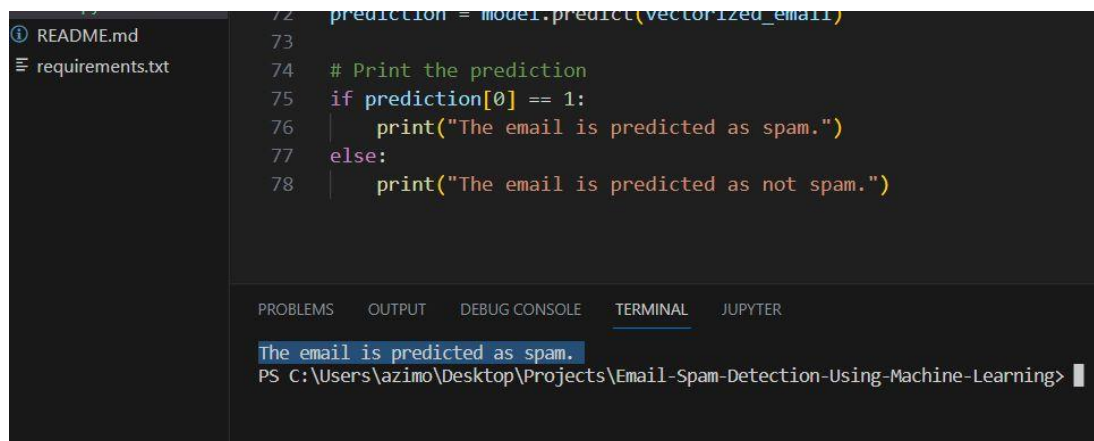


**Figure 3.** A visual representation of the amount of spam

Once the model is created and configured, it is stored in memory. Subsequently, the same model is loaded. Then, we use the test data to evaluate the established model [6] [7]. An example of the code for checking spam in the model is shown as follows (Figure 4).

```
# Example email to classify
email = "Subject: unbelievable new homes made easy  im wanting to show you this  homeowner  you h
"169 home loan at a 3 . 72 fixed rate .  this offer is being extended to you unconditionally and
"take advantage of this limited time opportunity  all we ask is that you visit our website and co
"look foward to hearing from you ,  dorcas pittman"
```

**Figure 4.** Test data

The resulting output is displayed below (Figure 5). If the submitted data is not spam, then the output will indicate that the data is a safe message and not spam [7] [10].

```
72    prediction = model.predict(vectorized_email)
README.md              73
requirements.txt       74    # Print the prediction
                       75    if prediction[0] == 1:
                       76        print("The email is predicted as spam.")
                       77    else:
                       78        print("The email is predicted as not spam.")


PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    JUPYTER
The email is predicted as spam.
PS C:\Users\azimo\Desktop\Projects\Email-Spam-Detection-Using-Machine-Learning>
```

**Figure 5.** Spam email detection

In the presented figures, in addition to the mentioned elements, the codes incorporating the prediction condition are also provided (Figure 6).

```
# Preprocess the email
def process_email(text):
    nopunc = [char for char in text if char not in string.punctuation]
    nopunc = ''.join(nopunc)
    clean = [word for word in nopunc.split() if word.lower() not in stopwords.words('english')]
    return clean

preprocessed_email = process_email(email)

# Vectorize the preprocessed email
vectorized_email = vectorizer.transform([' '.join(preprocessed_email)])

# Make prediction
prediction = model.predict(vectorized_email)

# Print the prediction
if prediction[0] == 1:
    print("The email is predicted as spam.")
else:
    print("The email is predicted as not spam.")
```

**Figure 6.** Data recognition through a vectorizer

After the model is constructed, it is essential to assess the extent to which it makes accurate and precise decisions. To achieve this, the following code example should be executed (Figure 7).

```
In [15]:   # Model saving
           dump(model, open("models/model.pkl", 'wb'))
```

```
In [16]:   # Model predictions on test set
           y_pred = model.predict(X_test)
```

```
In [17]:   # Model Evaluation | Accuracy
           accuracy = accuracy_score(y_test, y_pred)
           accuracy * 100
```

```
Out[17]:   99.20983318700614
```

**Figure 7.** Accuracy of the model

As evident from the figure, the accuracy metric of the created model is 99%. This implies that it can accurately classify 99 out of every 100 incoming emails.

Furthermore, it is possible to create a matrix that describes how accurately the model identified the normal and safe emails, how many spam emails were misclassified as non-spam, how many normal emails were identified as spam, and how many spam emails were correctly classified as spam among all the test data entered into this model [7] [8] [9]. The mentioned matrix is presented below (Figure 8).
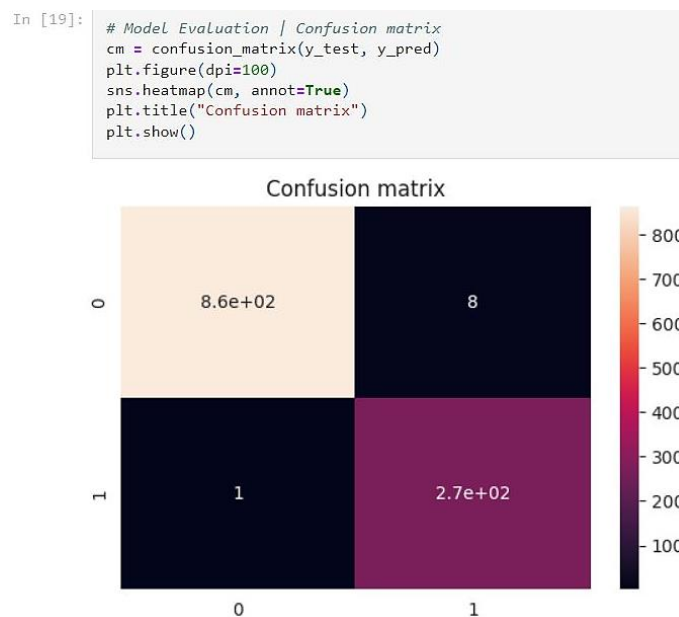
```
In [19]:   # Model Evaluation | Confusion matrix
           cm = confusion_matrix(y_test, y_pred)
           plt.figure(dpi=100)
           sns.heatmap(cm, annot=True)
           plt.title("Confusion matrix")
           plt.show()
```

**Figure 8.** Model evaluation

## III.    RESULTS

The protection of informational resources in companies holds significance not only for potential users but also for the companies themselves. When it comes to safeguarding databases and information sources, it is not merely achieved through software solutions. Hardware infrastructure establishment and utilization should also be considered. Corporations and organizations that ensure complete security safeguard their reputations and stand out sharply from their active competitors.

Each security model possesses distinct characteristics, and consideration should be given to their component composition when utilizing these models. Furthermore, to protect the company's assets from potential computer threats, security measures, such as maintaining and regularly updating signatures, should be diligently upheld. In the Python environment, spam protection detectors can effectively mitigate the problems encountered by companies. These detectors facilitate a more dynamic and robust defense against relevant attacks, thereby enhancing the overall security of the company.

REFERENCES

[1]    https://techvera.com/4-ways-to-protect-your-business-information-and-data/

[2]    https://www.endpointprotector.com/blog/5-ways-big-companies-protect-their-data/

[3]    https://www.fcc.gov/communications-business-opportunities/cybersecurity-small-businesses

[4]    Upton, David M.; Osborn, William J. Harvard Business School Cases. "Companies should understand where cybercrime thrives", Mar 01, 2022, p1-880. 880p.

[5]    Dobrygowski, Daniel. Harvard Business Review Digital Articles. "Why companies are forming alliances?", 9/11/2019, p2-5. 4p.

[6]    Watson, Rachel. Grand Rapids Business Journal. "Cyberattacks are gaining momentum", 5/2/2022, Vol. 40 Issue 9, p1-4. 4p.

[7]    Glassman, James K. Kiplinger's Personal Finance. "An urgent need for cybersecurity", Jun2022, Vol. 76 Issue 6, p39-40. 2p.

[8]    Aeilts Tony. FBI Law Enforcement Bulletin. "Defending against cybercrime and terrorism: A new role for universities." Volume: 74 Issue: 1 Dated: January 2005 Pages: 14-20

[9]    Peralta, Paola. Employee Benefit News. "Protect your company data from the Great Resignation", Mar/Apr2022, Vol. 36 Issue 2, pN.PAG-N.PAG. 1p.

[10]   Shaker, Bob. CFO. "Do You Even Know What Cyber Defense Is?", Jul/Aug2017, Vol. 33 Issue 6, p19-19. 2/3p.