



Department of Sociology

Master's Degree

Data Science

Computational Social Sciences Project

# **The importance of sentiment in scam filtering**

A study proposal on statistical and machine learning models  
that also include sentiment analysis

Student

Maurizio P. De Marchi (230654)

Academic Year 2021/2022

# INTRODUCTION

*“Phishing is the act of electronically and deceitfully contacting a person, with the aim of making the person electronically perform an act that is beneficial to the deceiver and harmful to the deceived. Phishing incidents typically occur by email”* [4]. In the current era, phishing is becoming more of a problem than ever. Phishing and scam attempts seem to have been growing [6], to the point that also companies famous for security awareness can fall victims [5]. Finding new ways to efficiently block them will be a key-point to avoid scammers’ social engineering acts and save web users from economical demise. Given the context, detecting and taking into consideration the social engineering attempt might be the key to improve filtering techniques. This proposal will focus on the emotional aspect of phishing social engineering, by attempting to incorporate it in predictive models.

## STATE OF THE ART

### Previous Studies

#### ***Phishing susceptibility***

The main reference for this section is *“A meta-analysis of field experiments on phishing susceptibility”* [4].

Phishing susceptibility is regarded as *“the probability that a recipient performs an action requested in a fraudulent message, is a widespread problem”*. Studying the psychological processes that lead a victim to succumb to a scam could be of great help to further augment existing email phishing prevention tools, however the literature about susceptibility appears to be scarce compared to the one about technical prevention.

In addition, no established model seems to exist to predict and understand phishing susceptibility.

Some studies seem to confirm existing general theories about deception and persuasion, where other evidences seem to contradict it.

Concluding this section, the results of the meta-analysis:

1. The characteristics of the person receiving the message (e.g.: personality, propensity to trust, etc.) seems to be of little significance.
2. The use of established deceptive tactics is confirmed to have an effect.
3. It seems that it is not important if the email comes from the right email server.
4. It was not explicitly tested if the content of the scam made a difference.
5. Circumstances are relevant, especially the timing: *“Data showed a lot of traffic due to the phishing emails just before closing hours and almost no traffic after closing hours”*. However, time of day, week and year were not controlled for.
6. Adaptation matters.

Other findings were more concerned about the methodological issues of the previous studies. In the hypothetical research project, dedicating some time to further research on this aspect could be of great importance, however it won’t be treated further here due to the scope of this assignment.

#### ***Emotionality***

Despite being already investigated by the previously cited meta-analysis, the two papers used for this section were deemed particularly useful due to the investigation of the textual content and anatomy of scams, thus shifting the perspective of the analysis.

According to *“You can trust me: a multimethod analysis of the Nigerian email scam”* [3] the concept of trust seems to be very common in the classic Nigerian email phishing template (also commonly referred as “419 scam” or “advance fee scam”).

Other important emotional manipulations that seem to be recurrent (according to “*The anatomy of written scam communications: An empirical analysis.*” [1]) are:

- Credibility, shown through care, concern and prioritizing the needs of others.
- Showing vulnerability, to give the impression of honesty.
- Compliments (social solidarity) and flattery, to lower the recipients’ guard.

## ***Spam filtering***

The main reference for this section is “*A Systematic Review on Spam Filtering Techniques based on Natural Language Processing Framework*”.

With respect to the spam filtering techniques, the following can be said (keeping in mind the Spam vs Ham dataset [7] context of the analysis):

- Despite not impacting greatly in susceptibility, the email address is regarded as a significant factor for detection.
- The email subject consistency tries to appeal to urgency.
- Given that redirecting the user to visit a link could be an important goal of the scammer, checking for link presence is deemed useful.
- Concerning Natural Language Processing (NLP), the analysis suggests the following:
  - Spam classification is done by identifying a list of repetitive words previously identified in scams
  - Other techniques such as Sentiment analysis and Topic Modeling were mentioned only with respect to spam comments. In addition, Polarity (positive or negative sentiment) seems to be the only aspect considered for spam, in the context of Sentiment Analysis.

## **Research Question**

Given the literature so far, it seems obvious that the emotional characteristic of text is an important part of scamming, however there is very little research on how emotionality of text contributes to the filtering process. Moreover, the studies claiming the importance of emotional manipulation (like [3]) have been conducted on a small corpus of spam texts. This context leads to the following research question: **what impact does text emotional connotation have in spam detection?** To answer the question, a big corpus of emails has been analyzed, including parameters such as sentiments, by an ensemble of statistical models. The whole process will be described in the following section.

## **METHODOLOGY**

### **Description of the dataset**

#### ***The making of***

Given the nature of the problem, various data gathering techniques have been used.

- To build the non-scam portion of the dataset, a portion of the famous enron email corpus [10] has been used. Given that having non-spam email was essential, only emails from the folders “\_sent\_mail” and “personal” were parsed into the dataset.
- To build the first part of the scam part of the dataset, all emails present on each page of 419scam.org [9] (2004 → 2022) have been scraped.
- To build the final, and most conspicuous part of the scam dataset, emails from untroubled.org [8] have been used. Given the great amount of data present, only emails from 1998, 2020 and 2022 have been parsed into the dataset.

After parsing and cleaning, the dataset was formed by 359.552 emails (289.337 scam, 70.215 not-scam).

## Data cleaning

Some cleaning operations have been performed:

- Emails with empty text, empty date or empty subject have been discarded.
- Non-UTF8 characters have been removed from the text.
- For sentiment analysis, email body texts have been converted to lowercase, stop-words have been removed, numbers and non-letters have been removed too. Finally, all words have been lemmatized.
- HTMLs tags have been removed.
- emails have been removed.
- common non-stop-words have been removed (enron, com, etc, ...);
- single characters have been removed

## Predictors

Given that text alone is not enough to build a model, predictors have been extracted from the fields *date*, *subject* and *text*.

Some however have been proven *problematic*:

- The **domain** of the email address is often treated with a whitelist/blacklist approach. Due to how the dataset was constructed, it is obvious that all emails coming from the enron domain are not-scam, so the domain could not be included as a predictor. To account for this, in the study one would have to find more non-scam emails. However, due to time constraints, in this proposal all emails have been assumed to have passed the domain check and needs to be identified after that.
- Email **urgency in the subject** is clearly an important aspect. Unfortunately, urgency detection in text is a hard problem of its own, so some surrogate predictors have been used.

Finally, the used predictors are:

- **Hour of the day** (approximated to the next after 45), *factor*;
- **Day of the week**, *factor*;
- **Day of the month**, *factor*;
- **Month**, *factor*;

- **Percentage of capital letter in the subject text** (*integer from 0 to 100*), in order to partially account from urgency.
- **Non-letter characters in subject** (spaces excluded), also to partially account for urgency. *integer*.
- **Length of the subject text**, *integer*.
- **Presence of links** in email text, *factor*.
- **Sentiment score in subject and text email**, *integer*. In particular:
  - anger;
  - anticipation;
  - disgust;
  - fear;
  - joy;
  - sadness;
  - surprise;
  - trust;
  - overall positiveness;
  - overall negativeness;

```
y = scam scam scam scam scam scam scam scam scam scam scam scam scam scam ... View
Month = may apr jun may mar mar may mar may feb apr apr mar apr jan mar may jan jul ... View
Day_number = 17 12 29 29 26 16 14 22 17 19 15 25 30 14 6 29 31 26 13 6 ... View
Day = mon fri thu thu tue mon mon sun wed thu tue fri sun mon mon thu fri fri mon ... View
Hour = 19 12 10 22 12 17 0 9 4 20 14 19 23 4 11 15 16 16 18 19 ... View
subject_capital_score = 7 16 1 4 11 14 9 10 56 7 3 25 15 50 6 56 50 16 6 20 ... View
subject_non_letter_score = 0 2 1 7 5 1 0 0 16 3 0 13 0 25 0 14 46 4 0 3 ... View
subject_length = 15 38 103 56 38 25 48 22 42 21 32 20 16 109 38 54 55 40 19 ... View
has_links = 1 0 0 0 1 1 1 0 1 1 1 1 1 1 0 0 0 1 1 ... View
Subject_anger = 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 ... View
Subject_anticipation = 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 ... View
Subject_disgust = 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 ... View
Subject_fear = 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 ... View
Subject_joy = 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 ... View
Subject_sadness = 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 ... View
Subject_surprise = 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 ... View
Subject_trust = 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 1 ... View
Subject_negative = 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 2 0 0 0 0 ... View
Subject_positive = 0 0 3 0 0 0 1 0 0 0 1 0 0 0 1 0 0 1 1 1 ... View
Text_anger = 2 0 2 42 17 0 0 0 0 0 0 0 1 0 0 0 0 8 0 0 ... View
Text_anticipation = 7 2 8 42 28 0 5 0 0 3 2 0 3 0 3 0 1 10 5 1 ... View
Text_disgust = 2 0 0 16 9 0 0 0 0 0 0 0 1 0 0 0 0 6 0 0 ... View
Text_fear = 1 0 1 47 22 0 1 0 0 0 0 0 0 0 0 0 0 5 0 0 ... View
Text_joy = 4 0 6 27 21 1 3 0 0 1 1 0 2 1 2 0 1 8 3 1 ... View
Text_sadness = 2 0 3 29 22 0 0 0 0 0 0 0 4 0 4 0 0 10 0 0 ... View
Text_surprise = 4 0 6 18 11 0 3 0 0 1 0 0 2 0 2 0 1 7 3 0 ... View
Text_trust = 7 2 9 69 27 1 7 0 0 1 1 0 5 1 5 0 1 11 6 1 ... View
Text_negative = 4 0 4 62 38 0 3 0 0 0 0 0 5 0 4 0 1 19 2 0 ... View
Text_positive = 10 2 18 110 46 1 11 0 0 2 1 0 9 1 9 0 2 13 9 ... View
```

Figure 1: predictors of the dataset

## Small exploratory analysis

As can be seen in figure 2 and 3, the most common words in scam emails are those regarding external resources, so one could expect the presence of links to be a very important predictor.



Figure 2: most common words on scam emails



Figure 3: most common words on non-scam emails

Another interesting thing to check preliminarily would be to look for a significant difference in sentiments. As can be seen in figure 4, the significant difference in mean between trust, positive, negative and overall sentiment in general would make one think of sentiment as a significant predictor.

```
Welch Two Sample t-test

data: df$Text_trust[df$y == "scam"] and df$Text_trust[df$y != "scam"]
t = 97.838, df = 268476, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 1.674837 3.825882
sample estimates:
mean of x mean of y
 7.928643 4.178884

Welch Two Sample t-test

data: df$Text_negative[df$y == "scam"] and df$Text_negative[df$y != "scam"]
t = 88.69, df = 243798, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 2.652339 2.772217
sample estimates:
mean of x mean of y
 5.764868 3.051798

Welch Two Sample t-test

data: df$Text_positive[df$y == "scam"] and df$Text_positive[df$y != "scam"]
t = 101.33, df = 258875, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 5.785718 6.813935
sample estimates:
mean of x mean of y
12.738835 6.838212

Welch Two Sample t-test

data: unlist(df[df$y == "scam", seq(from = 18, to = 37)]) and unlist(df[df$y != "scam", seq(from = 18, to = 37)])
t = 234.68, df = 4477966, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 1.126860 1.145841
sample estimates:
mean of x mean of y
 2.486328 1.269978
```

Figure 4: Significance t-test for average sentiments between scam and non-scam

## Analytical Strategy

The analytical strategy consists in training the models with the identified predictors. To grasp whether sentiments have an effect, three methods can be used:

1. For the parametric approaches, the use of the p-value.
2. For the non-parametric approaches, the use of variable importance.
3. Testing a model trained in different ways, in particular:
  - without sentiment predictors (control condition).
  - with sentiment predictors
  - only with predictors deemed significant by the p-value. Logistic regression has been used as the base to get those predictors, also given that variable importance had comparable results.

## Used models

1. **Logistic regression [18]**: due to the binary response, logistic regression was deemed as the most suitable linear model. This model identified as significant predictors (less than 0.01, see figure 2):
  1. some months
  2. the 22<sup>nd</sup> hour of the day
  3. non-letter scores
  4. the subject length
  5. if the email contains links
  6. text sadness
2. **Linear Discriminant Analysis [15]**
3. **Naïve Bayes [13]**
4. **K-nearest Neighbors [14]**
5. **Tree methods**, in particular:
  1. Simple classification trees [19] [see figure 6,7]
  2. Pruned trees [19] [see figure 8,9].
  3. Bagged Trees [16].
  4. Random Forests [16] [see figure 10].
  5. Boosted Trees [12]
6. **SVMs [13]**, in particular:
  1. Maximal Margin Classifier [see figure 11]
  2. Support Vector Classifier
  3. Support Vector Machine

- polynomial kernel
- sigmoid kernel
- radial kernel

## Figures

| Coefficients             | Estimate  | Std. Error | z value | Pr(> z )    |
|--------------------------|-----------|------------|---------|-------------|
| (Intercept)              | -3.958071 | 1.452836   | -2.719  | 0.00854 **  |
| Monthaug                 | -3.224805 | 1.375387   | -2.344  | 0.02007 **  |
| Monthdec                 | -3.262228 | 1.618908   | -2.020  | 0.04515 *   |
| Monthfeb                 | -8.162888 | 8.661588   | -0.942  | 0.34303     |
| Monthjan                 | -8.063188 | 8.692769   | -0.926  | 0.35243     |
| Monthjul                 | 2.771948  | 8.898992   | 0.312   | 0.75187     |
| Monthjun                 | 1.568133  | 8.735896   | 0.178   | 0.86408     |
| Monthmar                 | 8.335421  | 8.637783   | 0.965   | 0.33412     |
| Monthmay                 | 8.895815  | 8.661975   | 1.026   | 0.30656     |
| Monthnov                 | -1.866885 | 1.314886   | -1.420  | 0.15486     |
| Monthoct                 | -4.855585 | 1.576962   | -3.079  | 0.00288 **  |
| Monthsep                 | -2.283777 | 1.377876   | -1.657  | 0.09743     |
| Daymon                   | 1.338217  | 0.808877   | 1.643   | 0.09984     |
| Daytue                   | 0.817793  | 0.839768   | 0.974   | 0.33014     |
| Daythu                   | -0.374537 | 0.845513   | -0.443  | 0.65477     |
| Daythu                   | -0.388761 | 0.812827   | -0.478  | 0.63386     |
| Daywed                   | -0.908488 | 0.814588   | -1.103  | 0.26717     |
| Hour1                    | -1.284857 | 1.025624   | -1.253  | 0.21829     |
| Hour18                   | 0.144685  | 1.081897   | 0.134   | 0.89361     |
| Hour11                   | 1.897047  | 1.086382   | 1.747   | 0.08868     |
| Hour12                   | 0.886477  | 1.150416   | 0.771   | 0.44218     |
| Hour13                   | 2.927652  | 1.488274   | 1.968   | 0.04802 *   |
| Hour14                   | 1.296766  | 1.563585   | 0.829   | 0.40809     |
| Hour15                   | 2.413721  | 1.213268   | 1.989   | 0.04665 *   |
| Hour16                   | 2.922097  | 1.424171   | 2.052   | 0.04013 *   |
| Hour17                   | 1.874804  | 1.113128   | 1.687   | 0.09407     |
| Hour22                   | 3.444350  | 1.742095   | 1.980   | 0.04702 *   |
| Hour23                   | 4.394893  | 1.788526   | 2.458   | 0.00995 **  |
| Hour3                    | -0.898839 | 1.139670   | -0.792  | 0.43442     |
| Hour4                    | -0.377211 | 1.132840   | -0.331  | 0.74011     |
| Hour5                    | -0.888868 | 1.148854   | -0.776  | 0.43931     |
| Hour6                    | -0.255241 | 1.141232   | -0.224  | 0.82383     |
| Hour7                    | 0.207267  | 1.134271   | 0.183   | 0.85581     |
| Hour8                    | 0.170618  | 1.070728   | 0.159   | 0.87340     |
| Hour9                    | -0.756837 | 1.155660   | -0.657  | 0.51109     |
| subject_capital_score    | -0.813430 | 0.088784   | -9.150  | 0.00000 **  |
| subject_non_letter_score | 0.273749  | 0.056871   | 4.814   | 0.00000 **  |
| subject_length           | 0.037591  | 0.012654   | 3.018   | 0.00254 **  |
| has_links1               | 4.737163  | 0.522889   | 9.061   | < 2e-16 *** |
| subject_surprise         | -1.978610 | 0.901549   | -2.193  | 0.03468     |
| subject_trust            | -0.221585 | 0.442243   | -0.500  | 0.61696     |
| subject_negative         | 0.823748  | 0.489413   | 1.683   | 0.09314     |
| subject_positive         | 0.494850  | 0.346943   | 1.426   | 0.15434     |
| Text_sadness             | 0.816133  | 0.348798   | 2.340   | 0.02007 **  |
| Text_anger               | -0.987348 | 0.118281   | -8.346  | 0.00000 **  |
| Text_disgust             | 0.248858  | 0.248845   | 1.000   | 0.31728     |
| Text_fear                | -0.570835 | 0.187282   | -3.047  | 0.00254 **  |
| Text_joy                 | -0.188939 | 0.175316   | -1.076  | 0.28478     |
| Text_love                | 0.556878  | 0.197983   | 2.815   | 0.00257 **  |
| Text_surprise            | 0.118430  | 0.154681   | 0.766   | 0.44218     |
| Text_trust               | 0.868474  | 0.599115   | 1.449   | 0.14777     |
| Text_negative            | -0.192548 | 0.139816   | -1.379  | 0.16787     |
| Text_positive            | 0.918144  | 0.371118   | 2.474   | 0.01250 *   |

Figure 5: Logistic Regression p-values. Some non-significant have been cut due to space

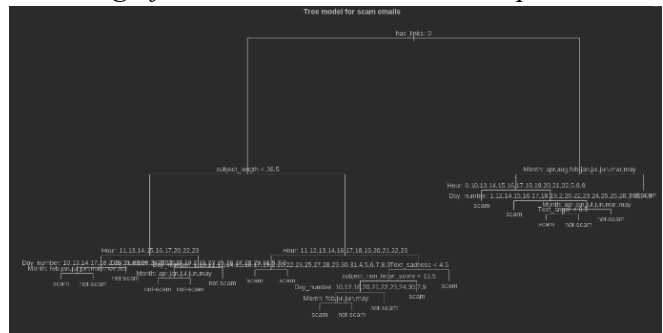


Figure 6: Simple classification tree

```

Classification tree:
tree(formula = y ~ ., data = df_tree)
Variables actually used in tree construction:
[1] "has_links"          "subject_length"      "Hour"
[4] "Day_number"         "Month"               "Text_sadness"
[7] "subject_non_letter_score" "Text_anger"

Number of terminal nodes: 19
Residual mean deviance: 0.2454 = 240.7 / 981
Misclassification error rate: 0.05 = 50 / 1000

```

Figure 7: predictors simple trees actually used

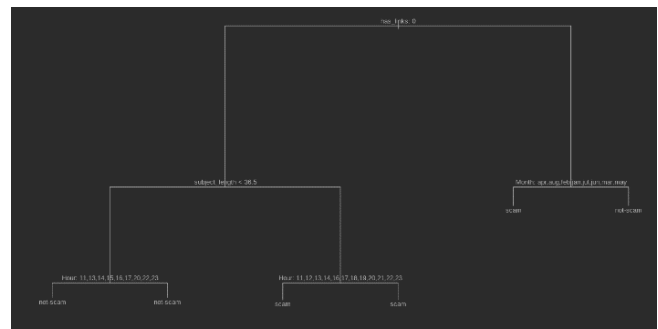


Figure 8: Pruned tree

```

Classification tree:
snip.tree(tree = model_tree, nodes = c(8L, 10L, 9L, 6L, 11L))
Variables actually used in tree construction:
[1] "has_links"      "subject_length" "Hour"           "Month"
Number of terminal nodes: 6
Residual mean deviance: 0.5256 = 522.5 / 994
Misclassification error rate: 0.104 = 104 / 1000

```

Figure 9: predictors pruned tree actually used

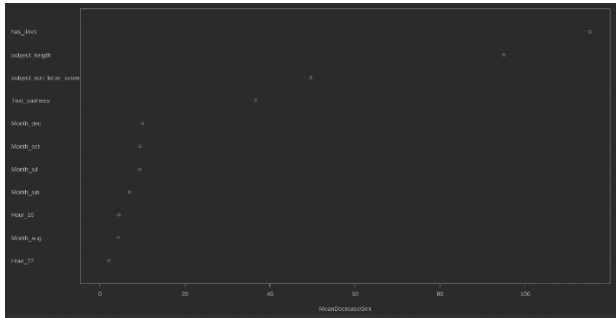


Figure 10: Random Forest variable importance

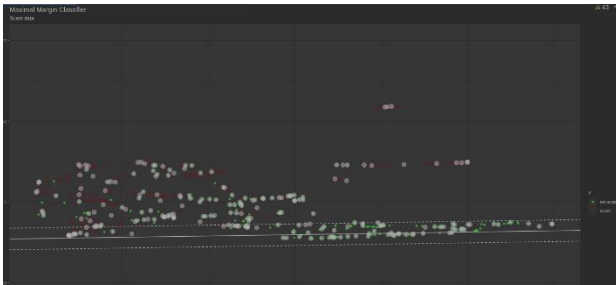


Figure 11: Maximal margin classifier (training)

## Example of analysis

Due to computational constraints, the example analysis has been conducted on a small sample of the dataset: 1000 emails for tuning and training, 1000 emails for testing (with a disproportion of 70/30 in favor of scam). The analysis was divided in model selection, where hyper-parameters were chosen through k-fold cross validation, and model assessment, where a test set was constructed by sampling new emails (never seen by the models) from the dataset. The process was done for each of the 3 settings (everything, no sentiment, only significant). Due to data uneven representation, the Area Under the Receiving operation curve (AUC) will be used as main metric reference, however other metrics (*Sensitivity*, *Specificity*, *Negative predicted value*, *Balanced Accuracy*, *Simple accuracy*, *Precision* and *F1 Score*) have also been computed for transparency.

As can be seen in tables 1 and 2, including sentiment did not significantly improve the predictive power of most models, with a difference of only ~1% of accuracy. Including only significant predictors also didn't improve the predictive power, in most cases it worsened it.

Table 1: All predictors (green=best, yellow=second, red=worst)

| Method                    | AUC   | Sensitivity | Specificity | Negative Predictive Value | Balanced Accuracy | Simple Accuracy | Precision | F1 score |
|---------------------------|-------|-------------|-------------|---------------------------|-------------------|-----------------|-----------|----------|
| Logistic Regression       | 95.0% | 92.3%       | 85.4%       | 81.7%                     | 88.8%             | 90.3%           | 94.0%     | 93.1%    |
| KNN                       | 91.9% | 90.2%       | 78.6%       | 77.0%                     | 84.4%             | 86.8%           | 91.0%     | 90.6%    |
| LDA                       | 95.5% | 94.3%       | 82.1%       | 87.0%                     | 88.2%             | 90.4%           | 91.9%     | 93.1%    |
| Naive Bayes               | 83.3% | 88.5%       | 42.4%       | 84.7%                     | 65.5%             | 60.9%           | 50.7%     | 64.5%    |
| Simple Tree               | 88.8% | 92.9%       | 74.6%       | 84.3%                     | 83.8%             | 86.7%           | 87.7%     | 90.2%    |
| Pruned Tree               | 90.7% | 89.9%       | 74.9%       | 76.7%                     | 82.4%             | 85.3%           | 89.0%     | 89.4%    |
| Bagging                   | 95.2% | 94.7%       | 76.8%       | 88.3%                     | 85.7%             | 88.5%           | 88.6%     | 91.5%    |
| Random Forest             | 97.9% | 95.9%       | 84.7%       | 90.7%                     | 90.3%             | 92.3%           | 93.0%     | 94.4%    |
| Boosting                  | 97.3% | 86.0%       | 93.7%       | 94.0%                     | 89.9%             | 91.4%           | 85.4%     | 85.7%    |
| Maximal Margin Classifier | 89.8% | 72.0%       | 81.6%       | 10.3%                     | 76.8%             | 72.4%           | 99.0%     | 83.4%    |
| Support Vector Classifier | 95.1% | 92.4%       | 84.8%       | 82.0%                     | 88.6%             | 90.2%           | 93.7%     | 93.0%    |
| SVM radial                | 92.3% | 83.5%       | 82.0%       | 56.3%                     | 82.8%             | 83.2%           | 94.7%     | 88.8%    |
| SVM Sigmoid               | 76.6% | 73.1%       | 49.6%       | 22.3%                     | 61.3%             | 69.9%           | 90.3%     | 80.8%    |
| SVM Polynomial            | 94.9% | 93.4%       | 84.7%       | 84.7%                     | 89.0%             | 90.8%           | 93.4%     | 93.4%    |

Table 2: No sentiment (green=best, yellow=second, red=worst)

| Method                    | AUC   | Sensitivity | Specificity | Negative Predictive Value | Balanced Accuracy | Simple Accuracy | Precision | F1    |
|---------------------------|-------|-------------|-------------|---------------------------|-------------------|-----------------|-----------|-------|
| Logistic Regression       | 94.7% | 92.4%       | 84.8%       | 82%                       | 88.6%             | 90.2%           | 93.7%     | 93%   |
| KNN                       | 92.7% | 93.3%       | 81.2%       | 84.7%                     | 87.2%             | 89.5%           | 91.6%     | 92.4% |
| LDA                       | 95.1% | 93.1%       | 82.1%       | 84%                       | 87.6%             | 89.7%           | 92.1%     | 92.6% |
| Naive Bayes               | 94.7% | 94.5%       | 80.7%       | 87.7%                     | 87.6%             | 90%             | 91%       | 92.7% |
| Simple Tree               | 88.7% | 90.6%       | 73.5%       | 78.7%                     | 82%               | 85.1%           | 87.9%     | 89.2% |
| Pruned Tree               | 88.8% | 89.9%       | 74.9%       | 76.7%                     | 82.4%             | 85.3%           | 89%       | 89.4% |
| Bagging                   | 94.2% | 89.2%       | 84.6%       | 73.3%                     | 86.9%             | 88%             | 94.3%     | 91.7% |
| Random Forest             | 96.9% | 89.7%       | 91%         | 74%                       | 90.3%             | 90%             | 96.9%     | 93.1% |
| Boosting                  | 96.2% | 86.7%       | 94%         | 94.3%                     | 90.3%             | 91.8%           | 86.1%     | 86.4% |
| Maximal Margin Classifier | 70.6% | 59.6%       | 15.4%       | 21.3%                     | 37.5%             | 41.2%           | 49.7%     | 54.2% |
| Support Vector Classifier | 95.1% | 92.5%       | 83.2%       | 82.3%                     | 87.8%             | 89.7%           | 92.9%     | 92.7% |
| SVM radial                | 95.7% | 92.2%       | 85.3%       | 81.3%                     | 88.7%             | 90.2%           | 94%       | 93.1% |



| Method         | AUC    | Sensitivity | Specificity | Negative Predictive Value | Balanced Accuracy | Simple Accuracy | Precision | F1    |
|----------------|--------|-------------|-------------|---------------------------|-------------------|-----------------|-----------|-------|
| SVM Sigmoid    | 86.8%  | 86.6%       | 78.1%       | 66.7%                     | 82.3%             | 84.4%           | 92%       | 89.2% |
| SVM Polynomial | 94.85% | 93.4%       | 80.7%       | 85%                       | 87.1%             | 89.4%           | 91.3%     | 92.3% |

Table 3: Only significant (green=best, yellow=second, red=worst)

| Method                    | AUC   | Sensitivity | Specificity | Negative Predictive Value | Balanced Accuracy | Simple Accuracy | Precision | F1    |
|---------------------------|-------|-------------|-------------|---------------------------|-------------------|-----------------|-----------|-------|
| Logistic Regression       | 94.5% | 92.4%       | 81.5%       | 82.3%                     | 87.0%             | 89.1%           | 92.0%     | 92.2% |
| KNN                       | 93.0% | 92.3%       | 82.6%       | 82.0%                     | 87.4%             | 89.4%           | 92.6%     | 92.4% |
| LDA                       | 94.5% | 94.5%       | 77.4%       | 88.0%                     | 86.0%             | 88.7%           | 89.0%     | 91.7% |
| Naive Bayes               | 92.4% | 93.5%       | 73.1%       | 86.0%                     | 83.3%             | 86.3%           | 86.4%     | 89.8% |
| Simple Tree               | 88.6% | 92.7%       | 77.2%       | 83.7%                     | 85.0%             | 87.7%           | 89.4%     | 91.1% |
| Pruned Tree               | 88.6% | 92.7%       | 77.2%       | 83.7%                     | 85.0%             | 87.7%           | 89.4%     | 91.1% |
| Bagging                   | 93.9% | 90.3%       | 79.7%       | 77.0%                     | 85.0%             | 87.2%           | 91.6%     | 90.9% |
| Random Forest             | 95.9% | 93.0%       | 80.3%       | 84.0%                     | 86.6%             | 89.0%           | 91.1%     | 92.1% |
| Boosting                  | 95.1% | 81.3%       | 92.4%       | 92.0%                     | 86.9%             | 89.1%           | 82.2%     | 81.7% |
| Maximal Margin Classifier | 74.4% | 70.0%       | NA          | NA                        | NA                | 70.0%           | 100.0%    | 82.4% |
| Support Vector Classifier | 94.2% | 93.0%       | 82.6%       | 83.7%                     | 87.8%             | 89.8%           | 92.4%     | 92.7% |
| SVM radial                | 93.7% | 93.1%       | 81.6%       | 84.0%                     | 87.3%             | 89.5%           | 91.9%     | 92.5% |
| SVM Sigmoid               | 81.3% | 83.0%       | 64.1%       | 59.0%                     | 73.6%             | 77.8%           | 85.9%     | 84.4% |
| SVM Polynomial            | 94.3% | 92.9%       | 81.8%       | 83.7%                     | 87.3%             | 89.5%           | 92.0%     | 92.5% |

## DISCUSSION

As seen so far, despite being an important aspect in the literature, text sentiment seems to be of little importance in a predictive context. This proposal's approach however has some weak points that needs to be fixed in the full study:

- Given that emails metadata is also important, more variety should be included. Taking data only from the enron dataset is a limitation from this point of view (e.g., only one domain for non-scam).
- Another limitation has been given by the available computational power: scam

emails were downsampled at the beginning to calculate sentiment in only few hours. Additionally, another downsample was done to tune models in a reasonable time. Finally, R [11] has been used as the main software to carry out computations, which can also be a limit because of the processor-centric logic (a package for GPU-based computations of SVMs has been tried, but R crashed without any error message).

- Temporal predictors that were not checked in previous studies seemed to be of little importance and cranks up computational times due to the high number of levels.

- A very big number of emails were not parsed due to the high variety in format. In the full study, more time should be dedicated to refining data parsing.
- As said in the methodology section, urgency detection is an aspect that was only partially accounted for.
- Different models and performance measures were used to avoid methodological biases. Tree based methods appears to be the most suited tool in the box for this problem, however there is still space for adding new models to test in the ensemble.
- P-value and variable importance were used to get the most significant predictors; however, the best-subset selection method could also be used with more computational power available.
- The library Syuzhet [17] was used to compute different kinds of sentiments. This is a double-edged sword: more sentiment variety to check is a clear advantage (e.g., text sadness was an important predictor), but makes dimensionality skyrocket, and so computational times. For reference, computing the sentiment of the dataset took ~ 2 hours.

In conclusion, sentiment analysis seems not to be so promising, but fixing the limitations could yield a different result.

All the code used for the assignment is available on Git Hub [20].

# SOURCES

## ***Papers and Meta-analysis***

- [1] Carter E. The anatomy of written scam communications: An empirical analysis. *Crime, Media, Culture*. 2015;11(2):89-103. doi:<10.1177/1741659015572310>
- [2] P. Garg and N. Girdhar, "A Systematic Review on Spam Filtering Techniques based on Natural Language Processing Framework," 2021 11th International Conference on Cloud Computing, Data Science & Engineering (Confluence), 2021, pp. 30-35, doi: 10.1109/Confluence51648.2021.9377042. <<https://ieeexplore.ieee.org/document/9377042>>
- [3] Rich, T. You can trust me: a multimethod analysis of the Nigerian email scam. *Secur J* **31**, 208–225 (2018). <<https://doi.org/10.1057/s41284-017-0095-0>>
- [4] T. Sommestad and H. Karlzén, "A meta-analysis of field experiments on phishing susceptibility", 2019 APWG Symposium on Electronic Crime Research (eCrime), 2019, pp. 1-14, doi: 10.1109/eCrime47957.2019.9037502. <<https://ieeexplore.ieee.org/document/9037502>>

## ***Articles***

- [5] Proofpoint phishing report  
<https://www.proofpoint.com/us/blog/security-awareness-training/2022-state-phish-explores-increasingly-active-threat-landscape>

- [6] Signal phishing attack:  
<https://support.signal.org/hc/en-us/articles/4850133017242-Twilio-Incident-What-Signal-Users-Need-to-Know->

## ***Cited datasets and data sources***

- [7] Spam vs Ham dataset  
<<https://www.kaggle.com/code/balakishan77/spam-or-ham-email-classification/data>>  
last visited: 04/08/2022
- [8] Spam dataset  
<<http://untroubled.org/spam/>>
- [9] Scraped emails  
<<https://www.419scam.org/emails/index.htm>>
- [10] Enron dataset  
<<https://www.cs.cmu.edu/~enron/>>

## ***R***

- [11] R software

**Libraries (NB: just the one cited in the report, not all libraries used for the project):**

- [12] caret  
<https://cran.r-project.org/web/packages/caret/index.html>
- [13] e1071  
<https://cran.r-project.org/web/packages/e1071/index.html>
- [14] kkn  
<https://cran.r-project.org/web/packages/kkn/index.html>

[15] MASS

<https://cran.r-project.org/web/packages/MASS/index.html>

[16] randomForest

<https://cran.r-project.org/web/packages/randomForest/index.html>

[17] syuzhet

<https://cran.r-project.org/web/packages/syuzhet/index.html>

[18] tidymodels

<https://cran.r-project.org/web/packages/tidymodels/index.html>

[19] tree

<https://cran.r-project.org/web/packages/tree/index.html>

## **Code**

[20] Git Hub repository

<https://github.com/DMMP0/Computational-Social-Sciences-project>