



UNIVERSIDAD DE CÓRDOBA  
ESCUELA POLITÉCNICA SUPERIOR DE CÓRDOBA

INGENIERÍA INFORMÁTICA  
ESPECIALIDAD: COMPUTACIÓN  
CUARTO CURSO. PRIMER CUATRIMESTRE

INTRODUCTION TO COMPUTATIONAL  
MODELS.

## Lab Assignment 4: Support Vector Machines (SVMs).

*Javier Herrero Porras*  
i72hepoj@uco.es

Academic Course 2020-2021  
Córdoba, December 1, 2020

# Contents

Figure index	ii
Table index	iii
1 Question 1	1
2 Question 2	2
3 Question 3	2
4 Question 4	3
5 Question 5	4
6 Question 6	5
7 Question 7	6
8 Question 8	7
9 Question 9	8
10 Question 10	8
11 Question 11	8
12 Question 12	9
13 Question 13	9
14 Question 14	10
15 Question 15	10
15.1 Emails wrong classified as Spam . . . . .	11
15.2 Emails wrong classified as Not Spam . . . . .	12
16 Question 16	15
17 Question 17	16

## List of Figures

1	Question 1 . . . . .	1
2	Question 3 . . . . .	3
3	Question 4 . . . . .	4
4	Question 5 . . . . .	5
5	Over-fitting and under-fitting . . . . .	5
6	Question 7 . . . . .	6
7	Question 7: Over-fitting and under-fitting . . . . .	7
8	Question 15: Confusion matrix . . . . .	11
9	Emails wrong classified as Spam (I) . . . . .	11
10	Emails wrong classified as Spam (II) . . . . .	12
11	Emails wrong classified as Spam (III) . . . . .	12
12	Emails wrong classified as Not Spam (I) . . . . .	13
13	Emails wrong classified as Not Spam (II) . . . . .	13
14	Emails wrong classified as Not Spam (III) . . . . .	13
15	Emails wrong classified as Not Spam (IV) . . . . .	14
16	Emails wrong classified as Not Spam (V) . . . . .	14
17	Emails wrong classified as Not Spam (VI) . . . . .	14
18	Emails wrong classified as Not Spam (VII)) . . . . .	15
19	Confusion matrix non-linear SVM . . . . .	17

## List of Tables

1	Accuracy dataset3.csv . . . . .	7
2	Accuracy and time Question 13 . . . . .	9
3	Confusion matrix . . . . .	16

# 1 Question 1

Open this script and explain its contents. You will see that the first dataset is used, and the SVM is graphically represented. Comment on what type of kernel is being used, and what are the training parameters. Explain the image generated, including all the different elements and colours.

In the script we train the SVM model with a linear kernel (that means that SVM model separate patterns with a straight line). Also we are using parameter  $C$ , which indicates the strength of the regularization, with the objective of controlling the error. If  $C$  is large that means that model accept less error and it will be a complex model. In our case  $C$  is large, so the model will accept less errors.

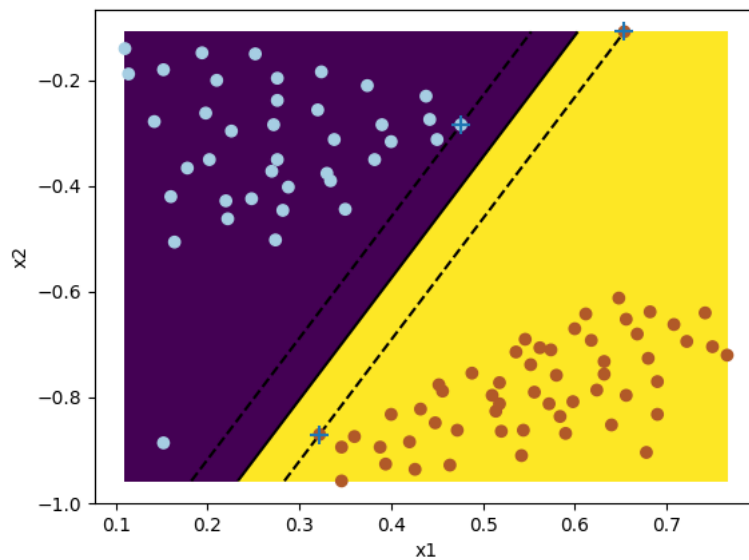


Figure 1: Question 1

In image 1 we can check the result of the model generated for the points. As we said before the model is very simple because we chose a linear kernel and  $C$  was large, so the model is able to accept some errors. We check that purple points are for example the blue ones and the other class is the red one. We can also verify the hyperplane obtained (continuous line) and the support vectors for this dataset (discontinuous lines).

## 2 Question 2

**Intuitively, which hyper plane do you think will make the least mistake in the task of separating the two kinds of points?**

In the image obtained, the simpler hyperplane which make the least mistake in separating the two kinds of points is the one which separates both classes without making mistakes. Also there is another option, which is the one that was made by SVM model in Question 1, which makes 2 mistakes in patterns.

## 3 Question 3

**Modify the script trying different values for  $C$ , specifically,  $C$  belonging to  $\{10e2, 10e1, 10e0, 10e1, 10e2, 10e3, 10e4\}$ . Observe what is happening, explaining why and select the most adequate value for  $C$ .**

In this case, I have changed the value of  $C$  by hand. In the next questions I will change the code to accept it by command line. With different values of  $C$  I have obtained these results:

- $C = 10e-2$ : the model is trying to maximize the margin, so it is accepting more errors. This lead to over-fitting, because the model is learning also noisy or wrong patterns.
- $C = 10e-1$ : the model has a small margin, so the errors made are less than model with  $C=10e-02$ . Now the model is more complex.
- $C = 10e0$ : the margin is smaller. The number of mistakes is reduced (10 now).
- $C = 10e1$ : the margin is smaller. The number of mistakes is reduced (5 now).
- $C = 10e2$ : the margin is smaller. The number of mistakes is reduced (3 now).
- $C = 10e3$ : the obtained model is the same as the previous value of  $C$ .

- $C = 10e4$ : the obtained model is the same as  $C = 10e2$  (mistakes = 3). The model can't make margin smaller. Maybe it's good to reconsider changing other model parameters to get a better model (kernel, gamma, etc).

For example, the model obtained with the value of  $C = 10e-0$  is shown in figure 2. In this case the most adequate value of  $C$  is  $10e2$ , because with higher values of  $C$  the model does not change.

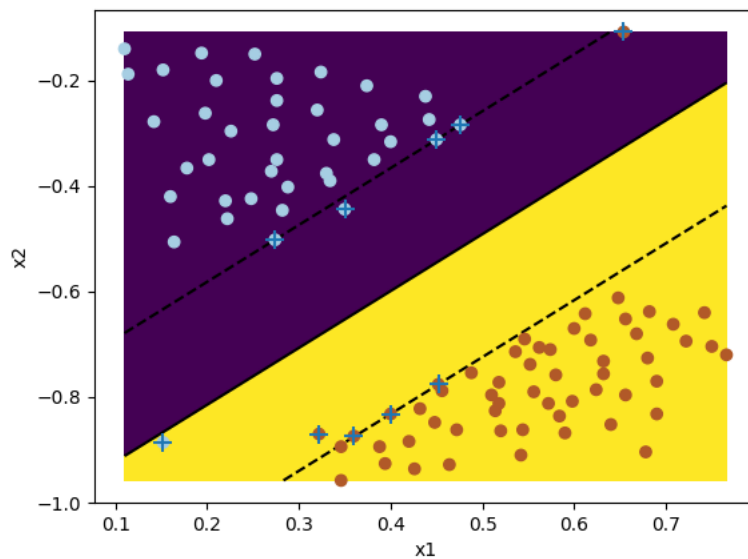


Figure 2: Question 3

## 4 Question 4

**Try running a linear SVM with the values for  $C$  used in the previous question. Do you get any satisfactory results in the sense that there are no errors in the training set? Why?**

In this case, data is not linearly separable (as we can see in figure 3), so that this model is not able to distinguish correctly between both classes. We will need kernel trick, to consider more dimensions in feature space and separate them. Then, we map the results into the input space, which is easier to compute in.

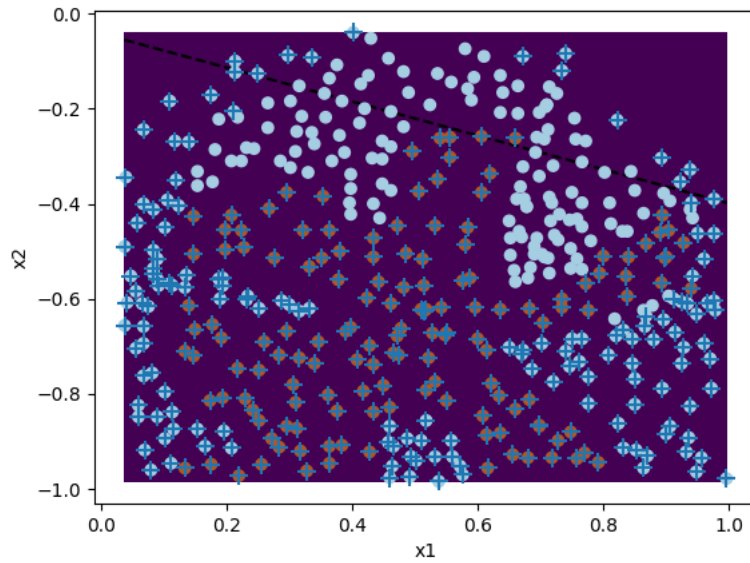


Figure 3: Question 4

## 5 Question 5

Propose a non-linear SVM configuration (using RBF or Gaussian kernel) that solves the problem. Try to use values for the parameters in the range  $C, \gamma$  belonging to  $\{10^{-2}, 10^{-1}, 10^0, 10^1, 10^2, 10^3, 10^4\}$ . What values have you considered for  $C$  and for  $\gamma$ ? Also, include an example of a parameter configuration that produces over-fitting and another that produces under-fitting.

For this question, I have chosen values for  $C = 1$  and  $\gamma = 10^2$ , obtaining the model shown in figure 4.



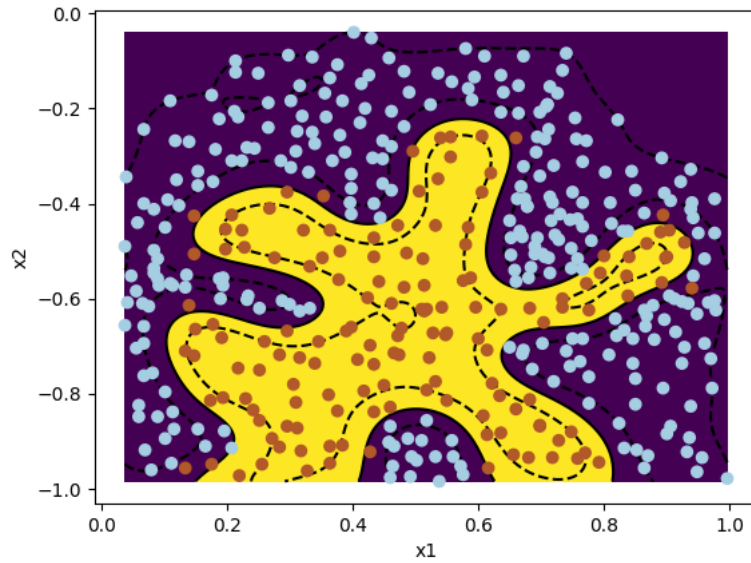


Figure 4: Question 5

One example of parameter configuration that produces over-fitting is  $C = 10^3$  and  $\gamma = 10^3$  and under-fitting is  $C = 1$  and  $\gamma = 1$ . The results are shown in figure 5.

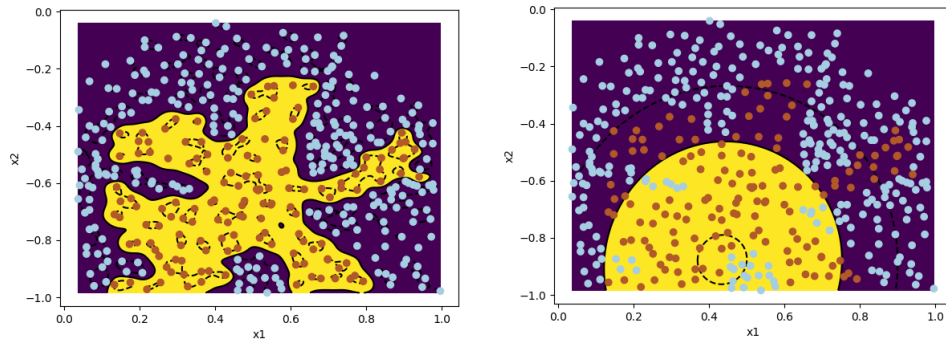


Figure 5: Over-fitting and under-fitting

## 6 Question 6

In this case, is the dataset linearly separable?. At first sight, do you detect points that are presumably 'outliers'? Why?

In this case, dataset is not linearly separable, because the points have the shape of two moons. Yes, I detect 2 red class points which are outliers, because its coordinates are equal to those ones which belong to the blue class.

## 7 Question 7

Run a SVM to classify the data, in order to obtain a result as close as possible to that of Figure 5. Adjust the parameters in the range  $C, \gamma$  belonging to  $\{10^{-2}, 10^{-1}, 10^0, 10^1, 10^2, 10^3, 10^4\}$ . Set the value of the optimal parameters. In addition, include an example of a parameter configuration that produces over-fitting and one that produces under-fitting.

The optimal parameters obtained are  $C = 10$  and  $\gamma = 1$  which produces the model shown in figure 6.

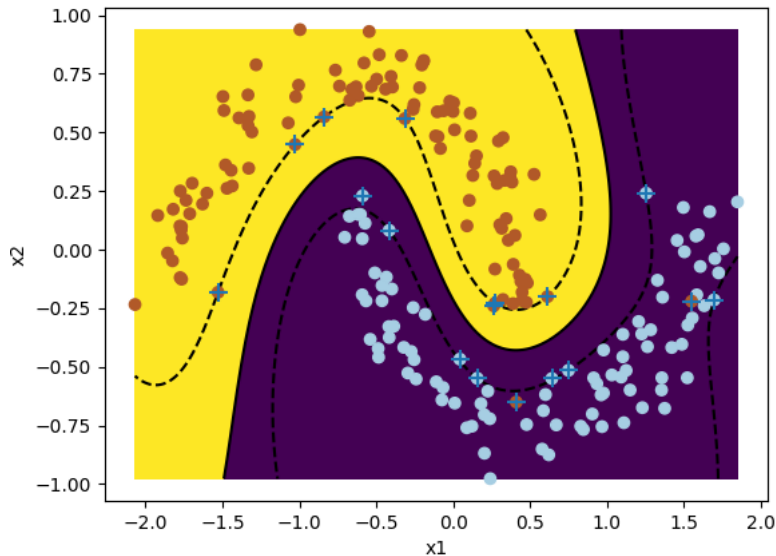


Figure 6: Question 7

One example of parameter configuration that produces over-fitting is  $C = 10^2$  and  $\gamma = 10^2$  and under-fitting is  $C = 10^{-1}$  and  $\gamma = 1$ . The results are shown in figure 7.

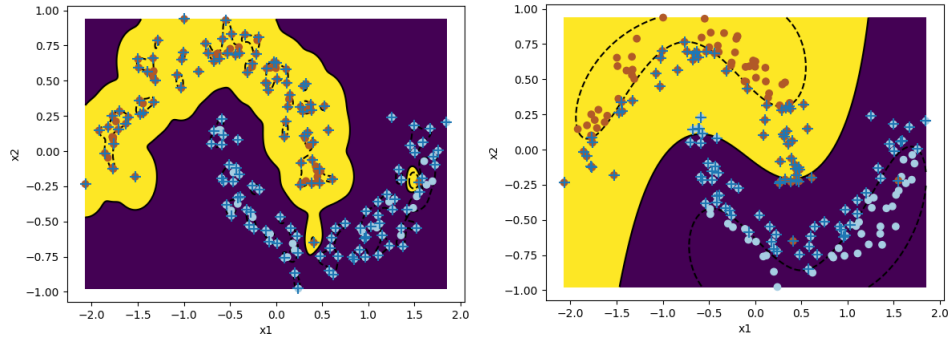


Figure 7: Question 7: Over-fitting and under-fitting

## 8 Question 8

We are going to reproduce this process in Python. Divide the synthetic dataset *dataset3.csv* into two stratified random subsets, with a 75% of patterns for training and a 25% of patterns for the test set. Make the complete training process (standardization, training and prediction), optimizing again the values of  $C$  and  $\gamma$ . Check the accuracy that is obtained for the test set. Repeat the process more than once to check that the results depend a lot on the seed used to make the partition.

In this question, I have obtained these results (with different executions, I choose one of the different parameter combinations which got a perfect accuracy). We can check that results depend on the random subsets taken by *train\_test\_split* function:

Execution	$C$	$\gamma$	Accuracy
1	100	1.0	1.0
2	100	1.0	0.96078
3	100	1.0	0.98039
4	100	1.0	1.0
5	100	1.0	1.0

Table 1: Accuracy dataset3.csv

## 9 Question 9

Extend the above code to perform the training of question 8 without the need to specify the values of  $C, \gamma$ . Compare the optimal values obtained for both parameters with those you obtained by hand. Extend the range of values to be explored, if you consider it necessary.

I remember that in question 8, I selected values  $C = 10^2$  and  $\gamma = 1$ . In this case, the optimal values obtained by *GridSearchCV* are  $C = 1$  and  $\gamma = 1$ , with an accuracy = 0.98039 (but with other execution the accuracy = 1, as happened in the previous exercise). These values simplify the model and also get a very good accuracy, so this method is working correctly.

## 10 Question 10

What drawbacks do you observe in adjusting the value of the parameters “by hand”, checking the accuracy in the test set (which was done in question 8)?

The first difference observed is that the grid search process is very fast in comparison with "by hand" method. Also, grid search is more secure than introducing parameter "by hand" because this is susceptible of errors made by the user.

## 11 Question 11

To be sure that you understand how the parameter search is performed, implement manually (without using *GridSearchCV*) the K-fold nested cross validation explained in this section. You may find useful the use of compression lists and the class *StratifiedKFold*. Compare the results with those you get using *GridSearchCV*.

The implementation made for this question is in file *Question11.py*. In this exercise, I have obtained an accuracy = 0.9868 (but it depends on the seed, because in other executions I have obtained an accuracy = 1). The optimal parameters are  $C = 1$  and  $\gamma = 1$ , so this demonstrates that this implementations is implemented properly.

## 12 Question 12

Use the script you developed in question 9 for the training of this database. Pay attention to the CCR value obtained for the generalization set and compare it to the one obtained in previous lab assignments. The process can take a long time. At the end, please, take a look to the optimal values obtained for the parameters.

In the execution, I have obtained an accuracy = 0.8866 and the optimal values for C and gamma are  $C = 100$  and  $\gamma = 0.001$ . This dataset is more complex than others studied before, so the model will be also more complex (due to the value of C) and the value of gamma (very low) implies that the model is not able to generalize new patterns which are far away from the training ones, because the model is constrained.

## 13 Question 13

Find out where the value of K for internal cross validation is specified and where the range of values used for the parameters C and  $\gamma$  are set. How could you reduce the computational time needed to carry out the experiment? Try to set K= 3, K= 5 and K= 10 and compare, using a table, the computational times obtained and the results for the test set in terms of CCR.

To set the value of K for the internal cross validation we have to use parameter *cv* in GridSearchCV. Also, for setting C and  $\gamma$  we have to set parameters C and gamma in GridSearchCV (the range of values is establish in *Cs* and *Gs*).

To reduce computational time needed for the experiment, we can change the parameter *n\_jobs* in GridSearchCV. In my case I set it to *n\_jobs* = -1, which indicates the maximum of *n\_jobs*.

K	Accuracy	Time
3	0.92	93.6937
5	0.8866	162.45
10	0.8866	374.32

Table 2: Accuracy and time Question 13

The result of the computational times and CCR is shown in table 2. We can see that accuracy is higher when we consider less  $K$ , that indicates us that the model works better with a lower number of training patterns (With  $K = 3$ , we consider 2/3 for training and 1/3 for validation), and also, the time increases when we increase  $K$  because it has to make more iterations and the process is also more complex.

## 14 Question 14

**A linear SVM model with the values  $C = \{10^{-2}, 10^{-1}, 10^0, 10^1\}$ . For this, use a script similar to the one used for question 9. Compare the results and establish the best configuration.**

The best configuration obtained with this script is a  $C = 0.01$  and  $\gamma = 0.001$  with a lineal SVM model. The accuracy obtained was 0.983. Although the dataset is more complex than others studied before, by the value of  $C$  we can say that the model will be simple and the value of gamma (very low) implies that the model is not able to generalize new patterns which are far away from the training ones, because the model is constrained.

## 15 Question 15

**For the best configuration, build the confusion matrix and establish the misclassified emails. Check the input variables for the emails incorrectly classified and find out the reason behind it. Note that for each pattern, when  $x_i$  is equal to 1 this means that the  $i$ -th word in the vocabulary appears, at least once, in the email.**

The implementation of the confusion matrix is done in file *Question15.py*. For the best configuration, we obtain this confusion matrix:

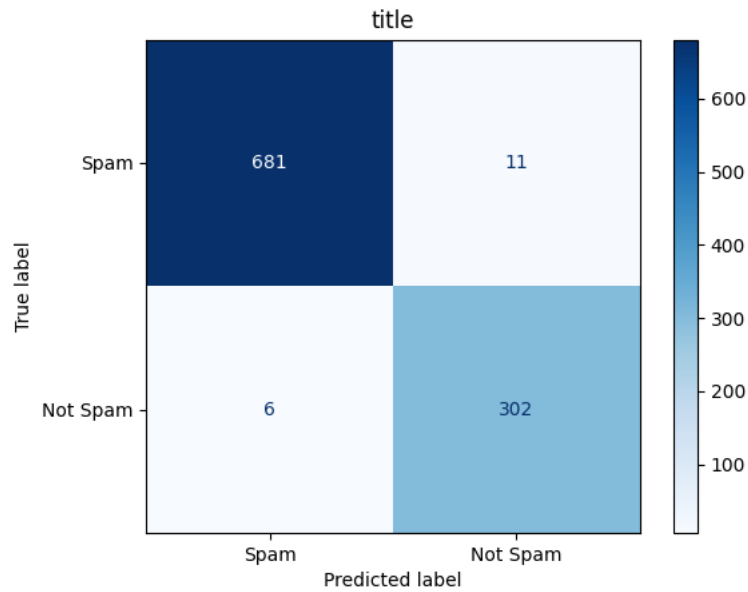


Figure 8: Question 15: Confusion matrix

In image 8 we check that there are 11 emails which are "spam" but the model classifies them as "not spam". Also there are 6 "not spam" emails that are classified as "spam". We are going to analyse these emails in the next subsection. We conclude that this model is able to classify correctly the "Not spam" emails with a high accuracy (98.05%), and the "Spam" emails (98.41%).

## 15.1 Emails wrong classified as Spam

The emails that are wrong classified as spam are 9, 90, 95, 150, 352 and 881. The content is shown in the following images:

```

9
[ 'emailaddr', 'for', 'group', 'httpaddr', 'inform', 'irish', 'linux', 'list', 'maintain', 'subscript', 'un', 'user' ]

90
[ 'aa', 'abil', 'abov', 'accord', 'account', 'action', 'add', 'administr', 'all', 'almost', 'alreadi', 'an', 'american', 'an', 'and', 'ar', 'ar', 'en', 'as', 'associ', 'at', 'averag', 'bar', 'be', 'been', 'best', 'better', 'black', 'board', 'bust', 'by', 'california', 'case', 'ceo', 'chan', 'c', 'chang', 'check', 'citizen', 'co', 'commun', 'consult', 'corpor', 'creat', 'current', 'dc', 'deliv', 'develop', 'differ', 'direct', 'down', 'educ', 'email', 'emailaddr', 'energ', 'engin', 'enjol', 'equal', 'everl', 'exclus', 'exist', 'far', 'feder', 'fight', 'for', 'former', 'fou', 'nd', 'from', 'geek', 'get', 'go', 'govern', 'ha', 'have', 'heaven', 'help', 'high', 'how', 'howev', 'httpaddr', 'id', 'immedi', 'in', 'indian', 'interest', 'is', 'it', 'let', 'list', 'lo', 'look', 'lot', 'mail', 'major', 'make', 'man', 'mean', 'measur', 'member', 'michael', 'monei', 'more', 'mostli', 'name', 'nation', 'net', 'new', 'no', 'non', 'not', 'number', 'numberf', 'or', 'on', 'onli', 'opportun', 'org', 'organ', 'ot', 'her', 'our', 'out', 'pacif', 'parti', 'past', 'peopl', 'perfect', 'perform', 'pictur', 'polit', 'potenti', 'present', 'presid', 'profession', 'progress', 'purpos', 'ratio', 'renov', 'repli', 'repres', 'return', 'seen', 'sf', 'show', 'sincer', 'situat', 'sourc', 'sponsor', 'state', 's', 'trategi', 'studi', 'support', 'tabl', 'take', 'templ', 'that', 'the', 'their', 'then', 'there', 'thi', 'thinkgeek', 'through', 'titl', 'to', 'toni', 'top', 'two', 'unfortun', 'univers', 'unsubscribe', 'us', 'vote', 'wa', 'washington', 'we', 'web', 'welcom', 'what', 'white', 'who', 'wi', 'll', 'win', 'with', 'woman', 'women', 'word', 'world', 'yourself' ]

```

Figure 9: Emails wrong classified as Spam (I)

```

95
['address', 'all', 'also', 'amp', 'an', 'and', 'ani', 'applic', 'ar', 'as', 'at', 'awai', 'base', 'bust', 'can', 'caus', 'charset', 'compani',
'comput', 'contact', 'content', 'databas', 'differ', 'emailaddr', 'encod', 'enhanc', 'etc', 'fact', 'find', 'for', 'format', 'from', 'futur',
'get', 'have', 'host', 'html', 'httpaddr', 'if', 'in', 'includ', 'inform', 'ireland', 'irish', 'is', 'iso', 'it', 'list', 'littl', 'local',
mai', 'mail', 'mani', 'messag', 'mime', 'month', 'more', 'multi', 'nbsp', 'nextpart', 'not', 'number', 'numberb', 'of', 'offer', 'onli', 'or',
'other', 'our', 'pack', 'part', 'per', 'plain', 'prefer', 'price', 'printabl', 'privaci', 'product', 'question', 'quot', 'receiv', 'remov',
'repli', 'requir', 'respect', 'respons', 'send', 'servic', 'shop', 'simpli', 'site', 'space', 'standard', 'state', 'support', 'sure', 'system',
'tax', 'text', 'the', 'these', 'thl', 'time', 'to', 'transfer', 'true', 'type', 'us', 'visit', 'we', 'will', 'with', 'without', 'would', 'ye',
', 'year', 'you', 'your']

150
['abov', 'account', 'after', 'against', 'all', 'also', 'an', 'amount', 'an', 'and', 'ani', 'anyth', 'ar', 'as', 'ask', 'assist', 'associ', 'av
ail', 'avoid', 'awai', 'ban', 'be', 'black', 'burn', 'but', 'by', 'can', 'care', 'choic', 'claim', 'close', 'come', 'commit', 'commun', 'compa
ni', 'confidentl', 'confirm', 'contact', 'content', 'convers', 'could', 'countri', 'current', 'death', 'deposit', 'develop', 'did', 'distribut
', 'do', 'document', 'down', 'drop', 'due', 'east', 'easill', 'email', 'emailaddr', 'enem', 'ensur', 'everyth', 'famill', 'far', 'father', 'f
ew', 'financ', 'follow', 'for', 'forward', 'found', 'free', 'friend', 'from', 'fund', 'futur', 'gain', 'give', 'god', 'govern', 'group', 'ha',
'hand', 'happen', 'have', 'he', 'help', 'here', 'hi', 'him', 'httpaddr', 'idea', 'immedi', 'in', 'inde', 'inform', 'instruct', 'interest', 'i
ntern', 'into', 'invest', 'irish', 'is', 'it', 'item', 'john', 'just', 'kept', 'kill', 'know', 'last', 'late', 'leav', 'left', 'linux', 'list',
'local', 'look', 'lot', 'mail', 'maintain', 'me', 'member', 'monet', 'move', 'mr', 'much', 'my', 'no', 'not', 'note', 'number', 'numberb',
'of', 'offic', 'on', 'or', 'order', 'our', 'out', 'over', 'own', 'pai', 'parti', 'person', 'phone', 'plan', 'pleas', 'polic', 'polit', 'present
', 'presid', 'press', 'privat', 're', 'reach', 'remain', 'repli', 'risk', 'robert', 'secur', 'share', 'should', 'so', 'son', 'sorri', 'steve',
'strong', 'subscript', 'sun', 'support', 'surpris', 'take', 'that', 'the', 'thet', 'their', 'then', 'there', 'therefor', 'thl', 'thing', 'tho
ugh', 'time', 'to', 'took', 'tool', 'total', 'transact', 'transfer', 'tri', 'two', 'un', 'under', 'upon', 'us', 'usdollarnumb', 'user', 'via',
'wa', 'wai', 'war', 'we', 'week', 'were', 'white', 'who', 'will', 'with', 'within', 'wonder', 'year', 'you', 'your', 'yourself']

```

Figure 10: Emails wrong classified as Spam (II)

```

352
['analyst', 'anoth', 'approach', 'ar', 'as', 'ascii', 'attent', 'back', 'be', 'bring', 'certain', 'charset', 'click', 'come', 'content', 'coul
d', 'dear', 'emerg', 'encod', 'everyon', 'for', 'format', 'gain', 'happi', 'have', 'here', 'html', 'huge', 'in', 'info', 'interest', 'is', 'is
o', 'itself', 'know', 'look', 'market', 'messag', 'mime', 'more', 'multi', 'nextpart', 'noth', 'number', 'numberbit', 'of', 'on', 'our', 'part
', 'plain', 'potenti', 'printabl', 'quot', 're', 'reader', 'record', 'rememb', 'sai', 'should', 'situat', 'sometim', 'speak', 'text', 'that',
'the', 'thet', 'their', 'thl', 'think', 'thought', 'to', 'track', 'transfer', 'type', 'us', 'ventur', 'we', 'when', 'will', 'with', 'you']

891
['about', 'abus', 'all', 'also', 'america', 'american', 'amount', 'an', 'and', 'ani', 'anoth', 'ar', 'argument', 'as', 'ask', 'attempt', 'atte
nt', 'author', 'awai', 'babi', 'base', 'basi', 'be', 'been', 'befor', 'behavior', 'believ', 'best', 'better', 'black', 'both', 'bui', 'by', 'c
ar', 'carri', 'caus', 'center', 'children', 'choic', 'choos', 'citizen', 'close', 'collect', 'con', 'commentari', 'commit', 'commun', 'confus',
', 'contain', 'cool', 'copi', 'could', 'coupl', 'cours', 'current', 'daili', 'date', 'did', 'direct', 'do', 'drive', 'drug', 'due', 'dynam', 'e
duc', 'email', 'emailaddr', 'even', 'ever', 'exist', 'expens', 'explain', 'ey', 'fact', 'femal', 'for', 'form', 'from', 'further', 'gener', 'g
o', 'govern', 'ha', 'have', 'heard', 'hi', 'higher', 'home', 'hous', 'how', 'howev', 'identifi', 'if', 'imag', 'import', 'in', 'individu', 'in
vest', 'is', 'it', 'know', 'known', 'lack', 'larg', 'learn', 'less', 'light', 'like', 'live', 'look', 'mai', 'make', 'maker', 'male', 'mani',
'mayb', 'me', 'men', 'might', 'miss', 'more', 'mother', 'move', 'my', 'name', 'natur', 'neg', 'new', 'no', 'not', 'noth', 'notic', 'now', 'num
ber', 'of', 'often', 'on', 'onli', 'or', 'over', 'own', 'pai', 'paid', 'part', 'particular', 'pass', 'peopl', 'perhap', 'place', 'pleas', 'poo
r', 'popul', 'prefer', 'project', 'purchas', 'rais', 'rather', 'refer', 'repres', 'request', 'reserv', 'respect', 'result', 'said', 'save', 's
elf', 'sell', 'shop', 'side', 'singl', 'sit', 'societi', 'soon', 'sort', 'stori', 'street', 'suffer', 'support', 'tax', 'than', 'that', 'the',
'thet', 'their', 'them', 'themselv', 'there', 'these', 'thl', 'through', 'to', 'todai', 'told', 'true', 'truth', 'util', 'visit', 'wa', 'warn
', 'we', 'wed', 'welcom', 'well', 'what', 'when', 'while', 'who', 'with', 'women', 'wonder', 'world', 'would', 'www', 'ye', 'you', 'young', 'y
our']

```

Figure 11: Emails wrong classified as Spam (III)

In general terms, it seems that our model detects spam on emails that have a few of words that begin by the same letter (for example a, t, etc). Also we can often find lots of pronouns, adverbs, prepositions, etc. In addition, the majority of emails have a lot of words.

## 15.2 Emails wrong classified as Not Spam

The emails that are wrong classified as spam are 21, 49, 73, 147, 208, 267, 328, 391, 526, 594 and 842. The content is shown in the following images:



```

21
['an', 'and', 'at', 'be', 'by', 'call', 'can', 'contact', 'email', 'emailaddr', 'geek', 'heaven', 'httpaddr', 'if', 'immedi', 'is', 'list', 'n
all', 'messag', 'need', 'net', 'not', 'number', 'of', 'office', 'our', 'out', 'repres', 'respond', 'return', 'sf', 'spamassassin', 'sponsor', '
start', 'talk', 'the', 'there', 'thi', 'thinkgeek', 'to', 'until', 'welcom', 'when', 'will', 'you', 'your']

49
['all', 'an', 'ar', 'ascii', 'at', 'bodi', 'charset', 'com', 'command', 'content', 'copyright', 'dai', 'distanc', 'eat', 'email', 'emailaddr',
'encod', 'english', 'find', 'friend', 'fun', 'go', 'he', 'here', 'hi', 'html', 'http', 'httpaddr', 'in', 'inc', 'intern', 'internet', 'is',
learn', 'look', 'messag', 'must', 'nbsp', 'number', 'numberbit', 'numberdumb', 'of', 'on', 'or', 'out', 'part', 'pictur', 'place', 'plain',
printabl', 'quot', 'reserv', 'right', 'see', 'send', 'six', 'text', 'that', 'the', 'thi', 'to', 'today', 'transfer', 'type', 'unsubscribe', 'us
', 'with', 'www', 'you']

```

Figure 12: Emails wrong classified as Not Spam (I)

```

73
['abl', 'about', 'abov', 'accept', 'access', 'account', 'act', 'addit', 'address', 'administr', 'advertis', 'advic', 'affect', 'affili', 'ag',
'again', 'against', 'agenc', 'agent', 'agre', 'air', 'all', 'allow', 'also', 'altern', 'an', 'among', 'an', 'and', 'anti', 'anoth', 'ar', 'are
a', 'as', 'associ', 'assum', 'at', 'attempt', 'author', 'auton', 'avail', 'award', 'be', 'becom', 'befor', 'begin', 'behalf', 'below', 'betwee
n', 'book', 'box', 'brand', 'broadcast', 'bug', 'busi', 'but', 'button', 'by', 'ca', 'california', 'can', 'cannot', 'capabl', 'case', 'cash',
'caus', 'center', 'centuri', 'chanc', 'check', 'choos', 'claim', 'clearli', 'click', 'co', 'collect', 'com', 'commit', 'commun', 'compani', 'c
omplet', 'comput', 'condit', 'confirm', 'connect', 'corpor', 'correspond', 'countri', 'court', 'dat', 'data', 'date', 'dear', 'decis', 'declar
', 'defin', 'deliv', 'depend', 'describ', 'direct', 'directli', 'director', 'do', 'document', 'doe', 'dollar numb', 'domain', 'doubl', 'draw',
'drive', 'dure', 'each', 'educ', 'either', 'electron', 'email', 'emailaddr', 'employe', 'enter', 'enterpris', 'entertain', 'entri', 'equal',
'equip', 'error', 'etc', 'event', 'except', 'expens', 'expir', 'fair', 'famili', 'feder', 'few', 'film', 'final', 'first', 'five', 'follow', 'f
or', 'foreign', 'form', 'free', 'from', 'full', 'fullt', 'game', 'gener', 'govern', 'grand', 'grant', 'greater', 'ground', 'ha', 'had', 'hardw
ar', 'hat', 'have', 'held', 'her', 'hi', 'him', 'hold', 'home', 'how', 'ident', 'if', 'immedi', 'improv', 'in', 'inc', 'includ', 'independ',
'indic', 'individu', 'inform', 'instal', 'institut', 'instruct', 'integ', 'internet', 'is', 'island', 'kind', 'late', 'later',
law', 'legal', 'like', 'limit', 'link', 'list', 'live', 'lo', 'local', 'locat', 'log', 'loss', 'lost', 'made', 'mai', 'mail', 'major', 'mani',
'manufactur', 'matter', 'meet', 'member', 'mention', 'messag', 'middl', 'militari', 'must', 'my', 'name', 'natur', 'near', 'necessari', 'nigh
t', 'no', 'non', 'not', 'number', 'numberd', 'numberth', 'oblig', 'occur', 'of', 'offer', 'office', 'offici', 'on', 'onc', 'onli', 'onlin', 'op
en', 'or', 'order', 'organ', 'origin', 'other', 'otherwis', 'out', 'over', 'oz', 'pack', 'page', 'paragraph', 'parent', 'part', 'parti', 'part
icip', 'per', 'period', 'permis', 'person', 'place', 'plan', 'pleas', 'point', 'polici', 'potenti', 'prevent', 'print', 'privaci', 'prize', 'p
roblem', 'process', 'product', 'promot', 'properti', 'provid', 'public', 'purchas', 'purpos', 'random', 'reason', 'receiv', 'regard', 'relat
', 'releas', 'repli', 'repres', 'request', 'requir', 'reserv', 'resid', 'respect', 'respons', 'result', 'retail', 'return', 'right', 'round',
rule', 'run', 'same', 'screen', 'script', 'second', 'secur', 'select', 'self', 'send', 'servic', 'set', 'sign', 'signatur', 'sincer', 'site',
'size', 'social', 'softwar', 'specif', 'specifi', 'sponsor', 'standard', 'state', 'subject', 'submit', 'such', 'take', 'taken', 'tax', 'techni
c', 'telephon', 'ten', 'term', 'termin', 'texa', 'than', 'thank', 'that', 'the', 'their', 'then', 'these', 'thi', 'thir', 'those', 'time', 't
n', 'to', 'togeth', 'traffic', 'transfer', 'travel', 'trip', 'under', 'unit', 'up', 'us', 'usa', 'valid', 'valu', 've', 'via', 'visit', 'voic
', 'wa', 'we', 'web', 'week', 'where', 'which', 'who', 'whole', 'whose', 'will', 'win', 'winner', 'with', 'within', 'without', 'won', 'write',
'www', 'you', 'your']

```

Figure 13: Emails wrong classified as Not Spam (II)

```

147
['about', 'accord', 'acquir', 'activ', 'actual', 'ad', 'addit', 'administr', 'advertis', 'advic', 'advic', 'after', 'again', 'ago', 'agre', 'a
lmost', 'alon', 'also', 'an', 'america', 'american', 'among', 'an', 'and', 'an', 'announc', 'annual', 'anoth', 'approach', 'ar', 'aren', 'aro
und', 'articl', 'artist', 'as', 'ask', 'associ', 'assum', 'at', 'august', 'averag', 'awai', 'back', 'bank', 'base', 'be', 'becaus', 'been', 'b
efor', 'below', 'better', 'big', 'billion', 'bit', 'board', 'both', 'bottom', 'brand', 'bring', 'bul', 'busi', 'but', 'by', 'cabl', 'call', 'c
an', 'capit', 'career', 'case', 'cash', 'caus', 'charg', 'check', 'chief', 'chip', 'chri', 'citi', 'civil', 'clearli', 'client', 'close', 'com
mit', 'commit', 'compani', 'complet', 'comput', 'conflict', 'connect', 'consid', 'construct', 'consult', 'consum', 'contact', 'contract', 'contro
l', 'copyright', 'corpor', 'cost', 'could', 'court', 'creativ', 'credit', 'cross', 'current', 'dal', 'david', 'deal', 'debt', 'decad', 'declin
', 'demand', 'describ', 'despit', 'detail', 'did', 'differ', 'director', 'discount', 'discov', 'discuss', 'do', 'doe', 'dollar', 'dollar numb',
'don', 'doubl', 'down', 'dure', 'earli', 'earlier', 'earn', 'easili', 'econom', 'econom', 'edg', 'electron', 'els', 'emailaddr', 'emerg', 'e
nd', 'entertain', 'environ', 'equip', 'establish', 'estim', 'even', 'ever', 'exampl', 'except', 'execut', 'exist', 'express', 'extens', 'face',
'fall', 'far', 'favorit', 'feder', 'few', 'file', 'financ', 'financi', 'find', 'firm', 'first', 'fix', 'flow', 'folk', 'for', 'form', 'forme
r', 'found', 'friend', 'from', 'fund', 'gener', 'global', 'go', 'goe', 'good', 'govern', 'grant', 'great', 'group', 'grow', 'guarante', 'gul',
'ha', 'had', 'half', 'hand', 'have', 'he', 'head', 'health', 'help', 'hi', 'high', 'highli', 'hold', 'home', 'hope', 'hour', 'how', 'httpaddr
', 'huge', 'human', 'if', 'in', 'includ', 'industri', 'inform', 'intit', 'instead', 'interest', 'internet', 'interview', 'into', 'invest', 'in
vestor', 'involv', 'is', 'issu', 'it', 'jame', 'jin', 'job', 'john', 'june', 'just', 'keep', 'king', 'know', 'knowledg', 'known', 'larg', 'lar
ger', 'last', 'late', 'later', 'latest', 'less', 'let', 'like', 'limit', 'line', 'littl', 'lo', 'loan', 'long', 'look', 'lose', 'lost', 'lower
', 'made', 'mai', 'mail', 'make', 'maker', 'manag', 'mani', 'march', 'mark', 'market', 'master', 'materi', 'media', 'men', 'million', 'monetl',
'month', 'more', 'most', 'mr', 'much', 'my', 'name', 'nation', 'nearli', 'necessari', 'need', 'network', 'new', 'newslett', 'no', 'nor', 'not
', 'note', 'now', 'number', 'numberb', 'octob', 'of', 'off', 'offer', 'office', 'often', 'old', 'on', 'onc', 'onlin', 'oper', 'opportun', 'or',
'organ', 'other', 'our', 'out', 'over', 'own', 'owner', 'pactif', 'paid', 'part', 'parti', 'partner', 'pattern', 'payment', 'peopl', 'percent',
'perfectli', 'perform', 'phone', 'pick', 'place', 'plat', 'plan', 'pleas', 'poor', 'pop', 'possibl', 'potenti', 'prefer', 'presid', 'pretti',
'price', 'privat', 'problem', 'profit', 'promis', 'prove', 'public', 'purchas', 'put', 'qualiti', 'quarter', 'question', 'quickli', 'rats',
're', 'reach', 'realiz', 'reason', 'recent', 'remain', 'replac', 'report', 'repres', 'result', 'retail', 'return', 'reveal', 'right', 'risk',
'robert', 'role', 'sai', 'said', 'sale', 'san', 'save', 'search', 'second', 'see', 'seed', 'seek', 'seen', 'self', 'sell', 'senior',
'sent', 'septemb', 'serv', 'servic', 'shape', 'share', 'should', 'show', 'sinc', 'sit', 'site', 'six', 'skeptic', 'skill', 'small', 'so', 'sold
', 'some', 'sometim', 'special', 'speed', 'spend', 'spot', 'stal', 'start', 'step', 'still', 'stock', 'structur', 'style', 'success', 'suffer',
'suit', 'system', 'take', 'taken', 'talk', 'team', 'technolog', 'texa', 'than', 'that', 'the', 'thei', 'their', 'them', 'then', 'there', 'th
ese', 'thi', 'think', 'thir', 'those', 'though', 'thought', 'three', 'through', 'ticket', 'time', 'to', 'toll', 'took', 'top', 'touch', 'trav
el', 'tri', 'trip', 'troubl', 'trust', 'turn', 'two', 'under', 'unless', 'unlik', 'until', 'up', 'upgrad', 'us', 'usual', 've', 'veri', 'visit
', 'vote', 'wa', 'want', 'war', 'washington', 'watch', 'we', 'web', 'week', 'well', 'were', 'west', 'what', 'when', 'where', 'which', 'who',
'whose', 'why', 'will', 'william', 'with', 'won', 'work', 'worth', 'would', 'write', 'ye', 'year', 'york', 'you', 'your']

```

Figure 14: Emails wrong classified as Not Spam (III)

```

208
['access', 'agent', 'allow', 'an', 'and', 'applic', 'ar', 'as', 'build', 'busi', 'by', 'commun', 'condit', 'configur', 'connect', 'consum', 'c
onveni', 'convers', 'current', 'custom', 'data', 'detail', 'develop', 'effici', 'end', 'everi', 'feel', 'few', 'fine', 'first', 'for', 'format
', 'four', 'function', 'get', 'great', 'ha', 'here', 'how', 'httpaddr', 'if', 'immedi', 'in', 'inc', 'includ', 'inform', 'instant', 'interact'
, 'interest', 'interfac', 'is', 'it', 'june', 'just', 'kind', 'launch', 'life', 'like', 'manag', 'mean', 'meet', 'messag', 'mike', 'million',
'more', 'natur', 'need', 'network', 'next', 'number', 'of', 'offer', 'on', 'oper', 'over', 'partner', 'patent', 'power', 'product', 'proven',
'provid', 'real', 'receiv', 'requir', 'sal', 'sampl', 'server', 'servic', 'softwar', 'solut', 'success', 'support', 'technolog', 'that', 'the'
, 'thi', 'time', 'to', 'track', 'type', 'unlimit', 'usag', 'user', 'virtual', 'visit', 'wa', 'what', 'with', 'wonder', 'year', 'you', 'your']

```

Figure 15: Emails wrong classified as Not Spam (IV)

```

328
['address', 'advertis', 'all', 'alwai', 'and', 'announc', 'been', 'below', 'but', 'close', 'cnet', 'com', 'comparison', 'custom', 'dal', 'down
load', 'each', 'emailaddr', 'faq', 'for', 'format', 'go', 'have', 'help', 'in', 'item', 'juli', 'let', 'mail', 'manag', 'million', 'more', 'my
', 'name', 'nbsp', 'new', 'number', 'of', 'on', 'our', 'own', 'price', 'product', 'provid', 'purchas', 'put', 'review', 'sale', 'save', 'servi
c', 'shop', 'so', 'subscrip', 'tech', 'technolog', 'that', 'the', 'there', 'thousand', 'to', 'today', 'unsubscribe', 'up', 'valu', 've', 'visi
t', 'you', 'your']

391
['about', 'abov', 'actual', 'after', 'almost', 'an', 'and', 'anoth', 'anyon', 'appear', 'ar', 'as', 'at', 'awar', 'ban', 'be', 'been', 'big',
'bodi', 'but', 'by', 'can', 'caus', 'compar', 'confer', 'consid', 'constitut', 'creat', 'current', 'degre', 'design', 'despit', 'detect', 'det
ermin', 'develop', 'di', 'discov', 'due', 'earth', 'energ', 'enter', 'equip', 'even', 'event', 'fact', 'first', 'flash', 'for', 'form', 'from
', 'ha', 'have', 'height', 'high', 'higher', 'hit', 'httpaddr', 'if', 'im', 'impact', 'in', 'increas', 'init', 'intern', 'into', 'is', 'it',
', 'juli', 'just', 'larg', 'light', 'look', 'low', 'lower', 'mat', 'mass', 'mayb', 'middl', 'more', 'network', 'no', 'number', 'occur', 'of', 'on',
', 'onli', 'or', 'part', 'peopl', 'planet', 'pre', 'pretti', 'probabl', 'public', 'reach', 'regist', 'rel', 'same', 'satellit', 'sever', 'stml
ar', 'size', 'some', 'start', 'techniqu', 'test', 'that', 'the', 'thel', 'these', 'thi', 'thu', 'to', 'us', 'usa', 've', 'wa', 'we', 'well', '
which', 'will', 'with', 'would', 'wrote']

526
['about', 'ad', 'advertis', 'after', 'am', 'an', 'and', 'as', 'at', 'been', 'below', 'bodi', 'but', 'by', 'click', 'contain', 'dat', 'direct',
'doc', 'email', 'emailaddr', 'except', 'extra', 'fals', 'final', 'geek', 'get', 'ha', 'have', 'heaven', 'here', 'httpaddr', 'is', 'it', 'list
', 'look', 'mail', 'mark', 'net', 'next', 'no', 'number', 'off', 'on', 'other', 'posit', 'razor', 'see', 'sender', 'sf', 'signatur', 'so', 'sp
an', 'sponsor', 'such', 'talk', 'the', 'thi', 'thinkgeek', 'to', 'turn', 'us', 'user', 'welcom', 'where', 'which', 'with', 'word']

```

Figure 16: Emails wrong classified as Not Spam (V)

```

584
['about', 'account', 'address', 'advertis', 'app', 'all', 'altern', 'american', 'among', 'an', 'and', 'anti', 'anim', 'announc', 'anyon', 'ar', 'as', 'at', 'auton', 'be', 'because', 'b
elow', 'best', 'by', 'can', 'check', 'con', 'come', 'comment', 'content', 'cool', 'cop', 'copyright', 'creat', 'current', 'cut', 'design', 'discov', 'discuss', 'distribut', 'do', 'd
ollar numb', 'don', 'dvd', 'each', 'earli', 'earn', 'easili', 'edit', 'email', 'end', 'enjoy', 'enter', 'entertain', 'evil', 'exampl', 'expert', 'face', 'favorit', 'featur', 'feel',
', 'film', 'fix', 'follow', 'for', 'forward', 'free', 'friend', 'from', 'front', 'full', 'futur', 'gain', 'game', 'get', 'give', 'great', 'group', 'ha', 'have', 'high', 'highli', 'home',
', 'hot', 'how', 'httpaddr', 'if', 'ill', 'in', 'inc', 'independ', 'inform', 'instruct', 'interest', 'internet', 'is', 'issu', 'it', 'judg', 'juli', 'just', 'keep', 'late', 'latest',
', 'launch', 'let', 'like', 'list', 'live', 'look', 'lot', 'nat', 'make', 'messag', 'million', 'more', 'now', 'near', 'new', 'newslett', 'next', 'not', 'now', 'number', 'of', 'on', 'ope
n', 'opt', 'or', 'origin', 'otherwise', 'out', 'over', 'own', 'page', 'pc', 'perfect', 'planet', 'player', 'pleas', 'popular', 'possibl', 'privat', 'protect', 'provid', 'qualiti', 're
ad', 'receiv', 'record', 'releas', 'remov', 'repli', 'request', 'reserv', 'result', 'retail', 'right', 'safe', 'sale', 'screen', 'section', 'seen', 'sell', 'send', 'sent', 'ser
v', 'set', 'ship', 'short', 'sign', 'site', 'skin', 'some', 'stal', 'stop', 'strong', 'style', 'subject', 'submit', 'subscrib', 'success', 'suggest', 'support', 'team', 'technolog',
', 'thank', 'that', 'the', 'thel', 'their', 'then', 'thi', 'thing', 'think', 'time', 'titl', 'tn', 'to', 'told', 'univers', 'unsubscribe', 'until', 'up', 'us', 'video', 'visit', 'wa', 'w
at', 'want', 'watch', 'we', 'week', 'weekli', 'welcom', 'what', 'while', 'who', 'wide', 'will', 'wish', 'with', 'world', 'worldwid', 'would', 'write', 'you', 'your']

842
['about', 'account', 'acquir', 'across', 'ad', 'address', 'advertis', 'affect', 'all', 'allow', 'almost', 'american', 'an', 'and', 'anti', 'anoth', 'anti', 'approach', 'ar', 'area',
', 'as', 'ask', 'asset', 'at', 'attack', 'avali', 'avoid', 'be', 'because', 'been', 'believ', 'best', 'better', 'bill', 'but', 'by', 'california', 'call', 'campaign', 'can', 'cent', 'ceo
', 'chanc', 'channel', 'chief', 'choic', 'choos', 'come', 'convert', 'commun', 'compan', 'congress', 'constitut', 'contact', 'content', 'could', 'countri', 'coverag', 'dal', 'decid
', 'defin', 'delet', 'differ', 'direct', 'directli', 'do', 'each', 'educ', 'effort', 'ensur', 'entir', 'equal', 'even', 'everi', 'expens', 'expert', 'face', 'fact', 'fair', 'featur',
', 'feder', 'few', 'field', 'firm', 'first', 'for', 'former', 'free', 'from', 'front', 'fund', 'futur', 'good', 'great', 'group', 'ha', 'had', 'have', 'he', 'hi', 'high', 'highli', 'him
', 'howev', 'httpaddr', 'if', 'in', 'individu', 'inform', 'instead', 'internet', 'is', 'it', 'just', 'know', 'lack', 'larg', 'law', 'least', 'less', 'let', 'level', 'light', 'like',
', 'link', 'list', 'lo', 'nat', 'mail', 'make', 'manag', 'nani', 'market', 'media', 'messag', 'method', 'night', 'mike', 'million', 'monet', 'month', 'more', 'most', 'must', 'nation',
', 'pag', 'no', 'not', 'now', 'number', 'of', 'off', 'offer', 'offic', 'on', 'onli', 'onlin', 'open', 'opt', 'or', 'other', 'our', 'out', 'paid', 'particularli', 'peopl', 'per', 'percent
', 'perfect', 'perform', 'person', 'phon', 'plai', 'point', 'polit', 'potenti', 'practic', 'prepar', 'presid', 'press', 'problem', 'provid', 'public', 'qualiti', 'radio', 'reach',
', 'read', 'receiv', 'recipi', 'refin', 'relev', 'remain', 'repli', 'repres', 'respons', 'seek', 'send', 'sender', 'sens', 'sent', 'she', 'shop', 'should', 'simpl', 'social', 'societi',
', 'softwar', 'some', 'span', 'spanassassin', 'specif', 'spot', 'staff', 'standard', 'state', 'such', 'support', 'take', 'taken', 'target', 'technolog', 'than', 'that', 'the', 'thel',
', 'their', 'then', 'thi', 'those', 'through', 'time', 'to', 'today', 'toward', 'tradit', 'tri', 'tv', 'unsolicit', 'unsubscribe', 'up', 'us', 'veri', 'vote', 'wa', 'wat', 'walk', 'want',
', 'war', 'we', 'what', 'when', 'while', 'who', 'whose', 'will', 'with', 'without', 'won', 'world', 'would', 'written', 'year']

```

Figure 17: Emails wrong classified as Not Spam (VI)



The test CCR obtained is 99.4% which is the best result obtained for the dataset. We can say that RBF network works better with this type of problem (but SVM has a lot of tuning parameters, and in this practice I have worked with a few of them, so maybe it would be a good idea to make more experiments with more tuning parameters). The confusion matrix obtained is:

	Spam	Not Spam
Spam	689	3
Not Spam	3	305

Table 3: Confusion matrix

From table 3 we can conclude that RBF model classify with a very high accuracy Spam e-mails (99,56%) and also is very good classifying "Not Spam" e-mails, because the accuracy is very high (99.02%).

## 17 Question 17

### **Train a non-linear SVM and compare the results obtained**

In this case, I have considered an 'rbf' kernel. After applying Grid Search, I have obtained the values  $C = 10$  and  $\gamma = 0.001$ , with a final accuracy of 0,968. The final confusion matrix is shown in figure 19. We conclude that this model is able to classify very well the "Spam" emails (accuracy = 99,56%) but it is not able to work correctly with predicting "Not spam" emails, because its accuracy decreases to 90.58%, so the model I propose in this Question would not be good for this dataset.

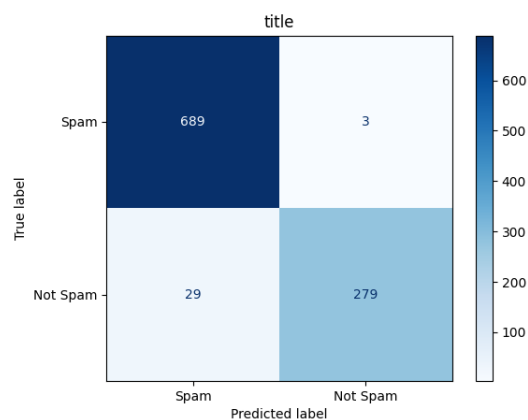


Figure 19: Confusion matrix non-linear SVM

We conclude that for "Spam" problem, a linear SVM model is going to classify better new parameters than for example an Gaussian SVM model, and for this problem the *kernel trick* is not useful at all because the data are lineally separable and if we use it, the accuracy decreases.