

Contents

Introduction	3
1 Statement of the Problem	5
1.1 Scope and Applications	5
1.2 Feature Selection	7
2 Basic Methods	8
2.1 Naive Bayes Classifier	8
2.2 k Nearest Neighbors Classifier	12
2.3 Artificial Neural Network Classifier	14
2.4 SVM Classifier	17
2.5 Boosting	19
3 Modified Methods	20
3.1 Non-Naive Bayesian Classifier	20
3.2 Combined Boosting	21
3.3 Adaptive Feature Selection	22
4 Performance Comparison in Practice	23
Conclusion	24
References	25

Introduction

The problem of spam has been around since the early days of e-mail services and has since evolved to numerous forms. Spam can vary by platform (e-mail, instant messaging, social network messaging), by intention (from advertisement to phishing and ransom), by distribution targets (from a single individual or general public), etc.

While the number of messaging domains is growing rapidly, the prime purpose of spam remains the same: to reach out to the recipient by any means or to trick him into acting on the message by disguising it as a legitimate message. This leads to a very clear practical definition of spam from user's perspective as any unsolicited message, whether it acknowledges its origin (e.g. harmless advertisement) or is crafted to gain the user's trust (e.g. phishing and ransom attacks).

Given the large number of messaging systems to which spam is applicable, the problem of detection of unwanted correspondence has evolved along with distribution methods and now includes additional metadata like attached media and hyperlinks. While they may seem to increase the complexity of spam detection, they actually give a basis for additional means of filtering which are often much easier to perform than the raw text classification.

For example, the vast majority of e-mail spam can be filtered out by simply checking the reputation of the sender's IP address and embedded hyperlinks. On the other hand, while it is possible to embed the text into images, their content can be easily retrieved by the use of optical character recognition and then appended to the main text.

While the message metadata can provide additional information, in practice text classification is inevitable. It is usually the final step in classification process as it tends to be most expensive. Since the metadata processing is different from

case to case, it is safe to say that the general problem of spam filtering comes down to the problem of binary text classification.

Therefore, in this paper we consider a number of common statistical methods of performing spam filtering on text messages and suggest a number of improvements to a number of them. Additionally, we analyze some of the common problems like feature vector selection. Finally, we provide a performance comparison of the detailed algorithms in practice.

Chapter 1

Statement of the Problem

1.1 Scope and Applications

A social networking service (or SNS) is a platform to build social networks or social relations among people who share similar interests, activities, backgrounds or real-life connections (definition from Wikipedia). The ultimate purpose of any social networking service is fast and efficient exchange of information, often with intention to present it to the largest audience possible.

The vast majority of most popular social network services rely on text messages as the main form of exchange of information. Because of their openness, social networks can be extremely useful in spreading malicious messages across wide audiences, both via private (addressed to a particular individual) and public messages. Conceptually social network spam is no different from e-mail spam as private messaging services of popular social networks are equivalent in their functionality to e-mail. Hence, we can generalize the problem of classifying spam messages for both these cases.

In general, the problem of spam detection depends heavily on the application's domain and can benefit from additional metadata available along with the text message. For example, in case of e-mails the mail header is the source of metadata. Modern spam filtering systems detect the vast majority of malicious mails by simply checking the sender's reputation before proceeding with analysis of the message body.

This, of course, applies to all messaging services. Maximum effectiveness can not be achieved without using all available data in addition to the message text.

However, in most cases text analysis is the second stage preceded by a domain-specific filter. Therefore, we can further focus on statistical classification of spam for text messages without specific constraints.

Let us denote the set of all messages by M , and let $S \subseteq M$ be the set of spam messages and $L = M \setminus S$ be the set of legit messages. The ultimate goal is to obtain a decision function $f : M \rightarrow \{S, L\}$ that would determine whether a given message $m \in M$ is spam ($m \in S$) or legitimate mail ($m \in L$).

We shall look for this function by training a number of machine learning algorithms on a set of already classified messages $\{(m_1, c_1), (m_2, c_2), \dots, (m_n, c_n)\}$, $m_i \in M, c_i \in \{S, L\}, 1 \leq i \leq n$. There are two aspects for the case of text messages: we have to extract features from text strings and we may have strict requirements for the precision of classifier.

1.2 Feature Selection

The entities we need to classify are text messages that are given in the form of strings. Raw strings are not convenient objects to handle in this case. Most machine learning algorithms can only classify numerical objects or otherwise require a distance metric or other measure of similarity between the objects.

Before proceeding with machine learning we have to convert all messages to numerical vectors called *feature vectors*, and then classify these vectors. The simplest example of a feature vector is the vector of the numbers of occurrences of certain words in a message.

Extraction of features usually means that some information from the original message is lost. On the other hand, the way feature vector is chosen is crucial for the performance of the filter. If the features are chosen so that there may exist a spam message and a legitimate mail with the same feature vector, then any machine learning algorithm will make mistakes no matter how good it is. A wise choice of features will make classification much easier while also fast. In most practical applications the most basic vector of word frequencies or its modification is used.

Note that at the stage of feature selection it is possible to include the features from the available metadata along with features from message text. In practice, however, it is much more important what features are chosen for classification than what classification algorithm is used.

Now let us consider those machine learning algorithms that require distance metric or scalar product to be defined on the set of messages. There does exist a suitable metric (edit distance), and there is a nice scalar product defined purely for strings [2], but the complexity of the calculation of these functions is sometimes too restrictive to use them in practice. So in this work we shall simply extract the feature vectors and use the distance/scalar product of these vectors.

Chapter 2

Basic Methods

2.1 Naive Bayes Classifier

Consider the simple case of text classification based on the presence or absence of just one word W . Suppose we know that the word W only occurs in spam messages. This gives us confidence that any message containing W is spam. This approach can be generalized to the probability of a message feature vector occurring in the message.

Suppose we have two classes L and S corresponding to legitimate and spam messages, and that there is a known probability distribution of feature vectors $P(x|c)$, $c \in \{L, S\}$. In general it is hard to define such distribution, but it is often possible to provide an approximation. What we need to obtain is the class that the given message belongs to, or the probability $P(c|x)$. This can be done using the Bayes' formula

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)} = \frac{P(x|c)P(c)}{P(x|L)P(L) + P(x|S)P(S)}.$$

where $P(x)$ is the a-priori probability of a message with feature vector x and $P(c)$ is the probability of class c , i.e. the probability that any given message belongs to c . Given the values $P(c)$ and $P(x|c)$ for $c \in \{L, S\}$ one can calculate the probability $P(c|x)$ which can then be used in a classification rule.

The most basic classification rule is to classify message to the category with bigger probability.

Definition 2.1.1. *Maximum a-posteriori probability (MAP) rule: if $P(S|x) > P(L|x)$ then classify x as spam, otherwise classify as legitimate message.*

The MAP rule can be transformed to

If $\frac{P(x|S)}{P(x|L)} > \frac{P(L)}{P(S)}$ then classify x as spam, otherwise as legitimate message.

The ratio $\frac{P(x|S)}{P(x|L)}$ is known as the *likelihood ratio* and is denoted as $\Lambda(x)$.

This approach can be too simplistic for certain applications. For example, in case of e-mail spam filtering, false positives (classifying legitimate message as spam) are usually much more unwanted than false negatives (classifying spam as legitimate message). The following generalization allows to take such restrictions into account.

Definition 2.1.2. A cost function $\mathcal{L}(c_1, c_2)$ denotes the cost of misclassifying a message of class c_1 as the one belonging to class c_2 .

It is natural to put $\mathcal{L}(L, L) = \mathcal{L}(S, S) = 0$, but in general it might not be the case.

Then we can express the expected risk of classifying a message x belonging to class c in the above terms.

Definition 2.1.3. The function $R(c|x) = \mathcal{L}(S, c)P(S|x) + \mathcal{L}(L, c)P(L|x)$, $x \in M$, $c \in \{L, S\}$ is called the risk function.

Now we can define a natural classification rule in terms of expected risk.

Definition 2.1.4. Bayes' classification rule: if $R(S|x) < R(L|x)$ then classify x as spam, otherwise as legitimate message [2].

It can be shown that Bayesian classifier f minimizes the average risk

$$R(f) = \int \mathcal{L}(c, f(x))dP(c, x) = P(L) \int \mathcal{L}(L, f(x))dP(x|L) + P(S) \int \mathcal{L}(S, f(x))dP(x|S)$$

so in this sense Bayesian classifier already is optimal [1].

Naturally, the loss of classifying the message correctly is zero, thus $\mathcal{L}(S, S) = \mathcal{L}(L, L) = 0$. Then the Bayes' classification rule can be rewritten as

If $\Lambda(x) > \lambda \frac{P(L)}{P(S)}$ classify as spam otherwise as legitimate message.

Here $\lambda = \frac{\mathcal{L}(L, S)}{\mathcal{L}(S, L)}$ is the additional parameter that specifies the risk of misclassifying legitimate messages as spam. As the value of λ increases, the classifier produces fewer false positives.

While the classification process is straightforward, the practical applications of Bayes's classifier are limited by our ability to approximate the a-priori probabilities $P(x|c)$ and $P(c)$, $c \in \{L, S\}$ from the training data. Therefore, while the Bayes's classifier is optimal in the sense of minimizing the loss of classification for given a-priori probabilities, the quality of spam detection depends on the feature selection and approximation of these probabilities.

$P(L)$ and $P(S)$ can be easily approximated by the ratio of legitimate and spam messages respectively. $P(x|c)$ is non-trivial and depends on the contents of selected feature vector. Consider the simplest case where the feature vector x_w is 1 if the message contains w and 0 otherwise. Then the probability $P(x_w = 1|S)$ can be approximated by the ratio of spam messages containing w to the ratio of all spam messages in a training set. This is sufficient to be used by the Bayes's classifier, so we can outline the training and selection process for this model.

Training process

1. Calculate probabilities $P(c)$, $P(x_w = 0|c)$, $P(x_w = 1|c)$, $c \in \{L, S\}$.
2. Using Bayes's formula calculate $P(c|x_w = 0)$ and $P(c|x_w = 1)$.
3. Calculate $\Lambda(x_w)$, $x_w = 0, 1$, calculate $\lambda \frac{P(L)}{P(S)}$ and store these values.

Classification process

1. Determine feature vector x_w for message m .
2. Retrieve the stored value $\Lambda(x_w)$.
3. Use Bayes's decision rule to determine class of the message.

Now we need to generalize this classifier to include more features than just the presence of a single word. The simplest way (and a very common one) is to choose a set of most common words w_1, w_2, \dots, w_n and define the feature vector $x = (x_1, x_2, \dots, x_n)$, $x_i = 1$ if the message contains w_i , $x_i = 0$ otherwise.

The problem with this approach is that it requires calculation and storing of all possible values of the feature vector, and there are 2^n such vectors, which is not feasible. A common way to remove this requirement is to assume that the individual components of the vector are independent [1]. This assumption is not formally correct, but in practice it is a good compromise between formal correctness and computational requirements. We will consider other options in later chapters.

Because of independence of features:

$$P(x|c) = \prod_{i=1}^n P(x_i|c), \quad \Lambda(x) = \prod_{i=1}^n \Lambda_i(x_i)$$

This classifier is known as Naive Bayesian Classifier due to assumption of independence of features. Training and classification are very simple computationally.

Training process

1. For all $w_i \in W$ calculate and store $\Lambda_i(x_i), x = 0, 1$.
2. Calculate and store $\lambda \frac{P(L)}{P(S)}$.

Classification process

1. Determine feature vector x for message m .
2. calculate $\Lambda(x)$ using the stored values $\Lambda_i(x_i)$.
3. Use Naive Bayes's decision rule to determine class of the message.

In terms of word selection for the feature vector, usually words that are too common or too rare are excluded. For simple cases when performance is not critical, all words from the training set can be used. In later chapters we will consider ways to select words with maximum mutual information.

Another benefit of naive Bayesian filter is that it is very easy to expand the feature vector to include additional available metadata. In case of e-mails, for example, it would be contents of e-mail headers. It is possible to include additional components either to the calculation of the a-priori probability of the vector or to combine the risk of Bayesian classifier with additional risk calculated from metadata when making a decision.

2.2 k Nearest Neighbors Classifier

k Nearest Neighbors (or k -NN) is a modification of the classical Nearest Neighbors algorithm. The idea behind this classifier is to first define a metric on feature vectors and then classify the message according to classes of k nearest messages in the training set. The metric is often chosen to be Euclidean, but Hamming or l_p can also be used for this purpose.

Training process

1. Store feature vectors of training messages in two sets L and S .

Classification process

1. For message with feature vector x determine k nearest neighbors from messages in the training set. If there are more legitimate messages among them, classify x as legitimate message, otherwise classify as spam.

Since the algorithm does not require any preprocessing of the training dataset, the training process is trivial. The classification process, however, requires calculation of distances to all messages in the training set, and for feature vectors of length m the most trivial implementations take $O(mn)$ time for set of n messages in case of Hamming or l_p metrics. Performing indexing on the training set can decrease the running closer to $O(n)$ [1]. However, if the size of the set increases over time, the algorithm might not be feasible in practice for certain applications.

The k -NN classifier is widely applicable in general classification problems, partially because it is one of so called *universally consistent* rules. Consider the training set s_n of n samples, and let us denote the k -NN classifier over that set as f_{s_n} . Similar to Bayesian classifier, we can determine the average risk $R(f_{s_n})$ of the classifier. The risk value is always greater than or equal than the Bayesian risk R_* (recall that Bayesian classifier is optimal in this sense), however for large values of n $R(f_{s_n})$ will be close to R_* .

Definition 2.2.1. *A classification rule is called consistent if the expectation of the average risk $E(R_n)$ converges to the optimal (Bayesian) risk R_* as n goes to infinity:*

$$E(R_n) \xrightarrow[n \rightarrow \infty]{} R_*$$

. The classification rule is called *universally consistent* if it is consistent for any distribution of (x, c) .

Theorem 2.2.1 (Stone, 1977). *If $k \rightarrow \infty$ and $\frac{k}{n} \rightarrow 0$, then k -NN classification rule is universally consistent.*

Consistency of k -NN rule allows to increase the quality by increasing the size of the training set. Stone theorem guarantees that as the size of the training set increases with constant value of k , the selection of messages for the training set does not matter. In addition to this, small values of k prevent quadratic complexity $O(n^2)$ when computing nearest neighbors.

Despite theoretical results, k -NN classifiers are performing worse than competition in practice in spam classification and are computationally expensive.

2.3 Artificial Neural Network Classifier

Artificial neural networks (ANN) is a family of models inspired by biological neural networks which are widely used in classification, regression and density estimation by approximating functions that can depend on a large number of inputs and are generally unknown. A neural network is a complex function that may be decomposed into smaller units called neurons and represented graphically as a network. Many functions fall under such criteria, however the most common kinds of neurons are perceptron and multilayer perceptron.

The perceptron produces a linear function of the feature vector $f(x) = w^T x + b$ where $f(x) > 0$ for vectors of one class and $f(x) < 0$ for vectors of another class. Here w is the vector of weights, or *bias*, $w = (w_1, w_2, \dots, w_n)$. This vector will be determined by the training process.

If we denote the classes by number -1 and +1, we can use $d(x) = \text{sign}(w^T x + b)$ as decision function. This allows us to represent the decision function graphically as a neuron with n inputs and a single output. A system of one perceptron is an example of the simplest neural network.

Suppose the feature vector is two-dimensional, $x \in \mathbb{R}^2$. Then we can represent these feature vectors as points on the plane. Then the decision function can be represented as a line dividing the plane in two parts, each corresponding to one of the classes. Similarly, the decision boundary for three-dimensional feature vectors is a plane, etc. In general, for n -dimensional feature vector the decision boundary is a n - dimensional hyperplane.

The learning process for a perceptron is iterative. The initial values of parameters (w_0, b_0) can be arbitrary, as they are updated on each iteration. On the k -th iteration of the algorithm a training sample (x, c) is chosen such that the current decision function does not classify it correctly, i.e. $\text{sign}(w_k^T x + b_k) \neq c$. Then the parameters (w_k, b_k) are then updated according to the rule:

$$w_{k+1} = w_k + cx$$

,

$$b_{k+1} = b_k + c$$

.

The algorithm terminates when a decision function that correctly classifies all training set is found. If the training set is linearly separable, the perceptron algorithm converges. It is known as Perceptron Convergence Theorem proven by Frank Rosenblatt in 1962 [4]. If, however, the set is not linearly separable, the algorithm will never converge. In this case it is possible to still use the perceptron, but the training loop needs to stop when the number of misclassification becomes small.

We can now outline the training and classification phases of the perceptron.

Training process

1. Initialize the values of w and b with random values or 0.
2. Find a sample from the training set (x, c) such that $\text{sign}(w^T x + b) \neq c$. If there are no such samples, terminate as the training is completed and all training samples are being classified correctly, else proceed to the next step.
3. Update (w, b) with new values $w := w + cx$, $b := b + c$ and go to the previous step.

Classification process

1. For message with feature vector x classify it as $\text{sign}(w^T x + b)$.

As mentioned before, perceptrons can be combined in multiple layers to form *multilayer perceptrons* which are non-linear classifiers. Neurons of the first layer which takes in the input parameters are called *input neurons*, similarly neurons of the last layer which provide the function result value are called *output neurons*. All layers between input and output are called *hidden layers*.

Each neuron in the networks is similar to a perceptron: for input vector $x = (x_1, x_2, \dots, x_n)$ it calculates output value by the formula

$$o = \phi\left(\sum_{i=1}^n w_i x_i + b\right)$$

where w_i and b are weights and bias of the neuron respectively, ϕ is a nonlinear function that approximates binary output of the perceptron. Most often $\frac{1}{1+e^{-ax}}$ or

$\tanh(x)$ are used as ϕ as they tend to give a good approximation while being mathematically convenient.

Like in the case of a single perceptron, training of a neural network is searching for the values of weights and biases for all neurons that minimize the output error. Let us denote $f(x)$ as the output of the neural network. Then for training samples $(x_i, c_i), 1 \leq i \leq k$ the training has to minimize the *total training error*

$$E(f) = \sum_{i=1}^k |f(x_i) - c_i|^2$$

An iterative algorithm can be used to perform this minimization. The most common one is the gradient descent which in case of neural networks is called *error backpropagation*. The theory of backwards propagation of errors is well developed and has many implementations in practice [3].

The main reason to use multiple layers of neurons is that the multilayer neural network is a non-linear classifier. As a result, they are applicable for tasks with training data that is not linearly separable, particularly when the number of features is relatively small. However, in case of spam detection with multiple words being used as features the data is often linearly separable, thus using neural network will have no noticeable benefits over a simple perceptron.

Performance of the neural network is proportional to the number of neurons. Thus, the large number of features directly impacts performance as it translates to increased number of input neurons and thus the complexity of the network in total. In practice the number of features would have to be more strictly limited than in case of a perceptron, which for spam detection means the trade-off between non-linear decision boundaries and the amount of information loss.

Because of the above reasons and due to the large number of parameters that require tuning, the multilayer perceptron is hard to use in practice for spam detection. It has been successfully used for that purpose [1], but it is not easily applicable in general case as it is hard to reconfigure. For the purposes of this paper we shall focus on a simple perceptron.

2.4 SVM Classifier

Support Vector Machines (SVM) is a class of widely used algorithms for classification and regression developed by V. Vapnik. The theoretical foundation of SVM is the Statistical Learning Theory that gives certain guarantees of performance. Let us consider classification problem with SVM for linearly separable data.

SVM works in a similar manner to perceptron in terms of finding a linear boundary that separates test data according to their classes. However, the purpose of SVM is not to find any of these boundaries if they exist, but to find the maximal margin separating hyperplane, for which the distance to the closest training sample is maximal.

Definition 2.4.1. Let $X = \{(x_i, c_i)\}$, $x_i \in \mathbb{R}^n$, $c_i \in \{-1, +1\}$ be the set of training samples. Suppose (w, b) is a separating hyperplane $\text{sign}(w^T x_i + b) = c_i$ for all $1 \leq i \leq k$. The margin m_i of a training sample (x_i, c_i) with respect to the separating hyperplane is the distance from x_i to the hyperplane

$$m_i = \frac{|w^T x_i + b|}{\|w\|}$$

The margin m of the separating hyperplane for training set X is the smallest margin of an instance in the training set

$$m = \min_i m_i$$

The maximal margin separating hyperplane for training set X is the separating hyperplane with maximal margin with respect to the training set [1].

Because the hyperplane given by parameters (x, b) is the same as the hyperplane given by parameters (kx, kb) , we can safely bound our search by only considering canonical hyperplanes for which $\min_i |w^T x_i + b| = 1$. It is possible to show that the optimal canonical hyperplane has minimal $\|w\|$, and that in order to find a canonical hyperplane it suffices to solve the minimization problem: minimize $\frac{1}{2}w^T w$ under conditions

$$c_i(w^T x_i + b) \geq 1, i = 1, 2, \dots, k$$

The problem may be transformed to a certain dual form: maximize

$$L_d(\alpha) = \sum_{i=1}^k \alpha_i - \frac{1}{2} \sum_{i,j=1}^k \alpha_i \alpha_j c_i c_j x_i^T x_j$$

with respect to dual variables $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_k)$, $\alpha_i \geq 0 \forall i$ and $\sum_{i=1}^k \alpha_i c_i = 0$.

This is a classical quadratic optimization problem. It mostly has a guaranteed unique solution, and there are efficient algorithms for finding this solution. Once we have found the solution α , the parameters (w_o, b_o) of the optimal hyperplane are determined as

$$w_o = \sum_{i=1}^k \alpha_i c_i x_i$$

,

$$b_o = \frac{1}{c_m} - w_o^T x_k$$

where m is an arbitrary index for which $\alpha_m \neq 0$.

It is more-or-less clear that the resulting hyperplane is completely defined by the training samples that are at minimal distance to it. These training samples are called support vectors and thus give the name to the method. It is possible to tune the amount of false positives produced by an SVM classifier, by using the so-called soft margin hyperplane and there are also lots of other modifications related to SVM learning. Recall the training and classification process.

Training process

1. Find α that solves the dual problem (maximizes L_d under named constraints).
2. Determine w and b for the optimal hyperplane and store the values.

Classification process

1. For message with feature vector x classify it as $\text{sign}(w^T x + b)$.

2.5 Boosting

Words, formulas.

Chapter 3

Modified Methods

3.1 Non-Naive Bayesian Classifier

Recall the naive Bayesian classifier described earlier. Let us ignore feature selection for now and instead only consider classification of feature vectors. Bayesian classifier is optimal in the sense of minimization of expected risk, however it is a common practice to assume that individual components of the feature vector are independent. Such assumption does compromise the optimality of the classifier, but also significantly simplifies implementation and improves performance. In this chapter we will analyze the possibility of an effective Bayesian classifier without this assumption and will develop an algorithm that allows to balance learning/classification speed with optimality guarantees.

3.2 Combined Boosting

Words, formulas.

3.3 Adaptive Feature Selection

Words, formulas.

Chapter 4

Performance Comparison in Practice

Words, tables, graphs, pictures, code.

Conclusion

Words.

References

- [1] Konstantin Tretyakov. Machine Learning Techniques in Spam Filtering. Institute of Computer Science, University of Tartu. Data Mining Problem-oriented Seminar, MTAT.03.177, May 2004, pp. 60-79.
- [2] V. Kecman. Learning and Soft Computing. 2001, The MIT Press.
- [3] S. Haykin. Neural Networks: A Comprehensive Foundation. 1998, Prentice Hall.
- [4] N. Cristianini, J. Shawe-Taylor. An Introduction to Support Vector Machines and other kernel-based learning methods. 2003, Cambridge University Press. <http://www.support-vector.net>
- [5] Xavier Carreras, Lluís Marquez. Boosting Trees for AntiSpam Email Filtering. TALP Research Center, LSI Department, Universitat Politècnica de Catalunya.