

Méthodologie d'évaluation des filtres anti-spam

José-Marcio Martins da Cruz
Mines ParisTech – Centre de Calcul et Systèmes d'information
60, bd Saint Michel
75006 - Paris
email: Jose-Marcio.Martins@mines-paristech.fr

Résumé

Nos boîtes aux lettres sont, depuis plus d'une dizaine d'années, pourries par les spams. Les administrateurs de messagerie utilisent, depuis une dizaine d'années, des outils divers, généralement des filtres distribués sous licence libre. On les met en place, assez souvent, grâce à des arguments qui nous semblent convaincants, ou par indications faites par des collègues, mais rarement leur efficacité est connue ou évaluée sérieusement. Parfois, agacés par les réclamations des utilisateurs, on met en place des nouveaux outils, juste dans l'espoir que la situation va s'améliorer. Le marché de la sécurité informatique, étant un marché très concurrentiel, les fournisseurs ne renseignent pas toujours les vraies indices d'efficacité de leurs produits.

« Tout le monde le sait » : l'efficacité d'un filtre anti-spam se mesure par ses taux de faux positifs et faux négatifs. Mais que veulent réellement dire ces chiffres ? Ces indicateurs sont ils suffisants et adaptés à l'évaluation d'un filtre anti-spam ? Comment évaluer sérieusement un filtre anti-spam ?

Cet article présente une méthode d'évaluation de filtres anti-spam, prioritairement adaptée pour l'évaluation de filtres de contenu, mais qui peut être adaptée pour les filtres protocolaires. Dans une première partie nous passons en revue les particularités du filtrage de spam, les points importants à tenir compte dans une évaluation ainsi que les indicateurs d'efficacité de classement usuels et pertinents pour le cas de filtrage de spam.

Mots clefs

Filtrage de spam, évaluation, filtrage en ligne

1 Introduction

Dans la deuxième moitié des années 90, la problématique du filtrage de spam, par son contenu, a attiré l'attention de trois domaines de recherche qui n'étaient pas directement concernés par la messagerie électronique : la Recherche Documentaire (IR), la Fouille de Données (DM) et l'Apprentissage Automatique (ML)¹. Le premier résultat de recherche publié d'utilisation d'un filtre statistique pour filtrer le spam date de 1998, par Sahami et al [1], utilisant un classificateur bayésien naïf. En 2002, Paul Graham a publié un billet dans son blog [2] qui a déclenché un foisonnement de filtres statistiques, dits bayésiens, distribués sous licence libre. Plusieurs centaines de communications scientifiques ont été publiées depuis, et nombreux logiciels de filtrage sont disponibles (des produits commerciaux ou libres).

On est souvent confronté à des affirmations du genre : « ce filtre anti-spam est capable de détecter 99 % des spams sans aucun faux positif ». Cette information semble intéressante mais, sauf à avoir plus de précisions, elle n'est pas vraiment utile. Quel administrateur de messagerie ne s'est jamais posé la question : quelle est l'efficacité de mon filtre anti-spam ? Si j'essaye telle ou telle modification dans la configuration de mon filtre, comment vérifier si cela améliore l'efficacité ? Comment comparer deux filtres anti-spam ? Ou tout simplement, de combien j'améliore l'efficacité du service de filtrage si j'ajoute un filtre en plus ?

Les chercheurs et les développeurs ont besoin d'évaluer les filtres pour mettre au point et sélectionner des méthodes de filtrage. Les gestionnaires de services de messagerie ont le même besoin d'abord pour décider lequel convient le mieux à son organisation et après pendant la phase de production, pour détecter toute déviation par rapport à l'efficacité initiale.

Le marché de la sécurité informatique, en général, et de la messagerie, en particulier, est très concurrentiel et on ne doute pas que les fournisseurs soient réticents à ce que leurs produits puissent être évalués et comparés avec d'autres.

Vaderetro indique, dans la plaquette de présentation de Mail-Cube [3], une « *Détection Importante des spams (95 %)* » mais ne précise ni les conditions d'évaluation ni à quel taux de faux positifs. Il n'y a pas plus d'information dans le manuel d'utilisateur. De même, il est indiqué que « *plus de 100 messages à la seconde peuvent être analysés sur une machine équipée d'un Pentium 4 cadencé à 1,9 GHz* », mais pas de précision sur la taille des messages, paramètre souvent crucial pour la capacité de traitement, ou même si le filtre partage une machine avec le MTA ou s'il est seul.

¹Ces domaines sont connus dans la littérature en langue anglaise par les noms « Information Retrieval », « Data Mining » et « Machine Learning ».

La plaquette commerciale de Ironport [4], se limite à dire que le boîtier est capable de « *refuser jusqu'à 80 % du spam dans la phase de connexion, sans aucun risque de faux positifs* ». Il n'y a pas mention sur l'efficacité du filtrage de contenu. Dans les messages traités par le filtre de contenu, des entêtes sont ajoutés mais les informations sont codées, ce qui rend difficile (mais pas impossible) toute évaluation faite par des tiers, utilisant les valeurs de score.

Lors de la « Spam Conférence » du MIT en 2004, Bill Yerazunis a fait une présentation avec le titre « *The Spam-Filtering Accuracy Plateau at 99.9 % Accuracy and How to Get Past It* » [5]. Un an après, lors de la conférence TREC – Spam Track 2005 [6], quatre configurations de son filtre ont été testées contre quatre corpus de messages : les résultats de précision globale² ont varié entre 87,95 % et 99,66 % et n'ont donc pas atteint le niveau annoncé par l'auteur. L'utilisation d'exemples d'apprentissage dans l'ensemble à tester et la simulation sans tenir compte de l'ordre chronologique des messages ne serait valable que si le spam était un processus stationnaire : ce sont les erreurs commises par Yerazunis.

La différence entre les mesures effectuées peut même donner naissance à des situations conflictuelles. Gordon Cormack (professeur à l'Université de Waterloo) a comparé l'efficacité de 6 filtres anti-spam open source (Bogofilter, SpamAssassin, Dspam, CRM114, SpamProbe et SpamBayes) [7]. Zdziarski, auteur de Dspam, non content de voir son filtre classé en avant dernière place (derrière SpamAssassin et Bogofilter) a démarré une polémique. Pourtant, le protocole d'évaluation utilisé par Cormack est largement plébiscité par la communauté de la recherche, alors que Zdziarski [8] n'évalue que la précision globale (accuracy), paramètre qui n'est pas pertinent pour une application de classement tel le filtrage de spam [9], avec un protocole d'évaluation moins justifié que celui de Cormack.

Les caractéristiques intrinsèques des filtres (par exemple, le type de classificateur et les caractéristiques d'extraction des attributs) font que certains filtres donnent des meilleurs résultats que d'autres. Ajoutons que les facteurs externes (par exemple, les jeux de données, le protocole d'apprentissage, l'environnement linguistique) font que le classement d'un filtre peut varier d'une évaluation à une autre.

Il est donc nécessaire, non seulement d'évaluer les filtres dans des conditions le plus proche possible des conditions réelles d'utilisation, mais aussi de les préciser. Il est parfois aussi intéressant d'évaluer la sensibilité du filtre à des changements limités des conditions d'évaluation avec, par exemple, ajout de bruit.

Le but de ce papier est la présentation des difficultés de l'évaluation des filtres anti-spam, des mesures pertinentes et de présenter les idées de méthodologie utilisée dans TREC Spam Track, méthodologie qui semble être actuellement la plus acceptée par la communauté de la recherche.

Il y a deux catégories de filtres anti-spam³ :

- **filtres protocolaires ou comportementaux** - ce sont des filtres dont le fonctionnement est basé sur des paramètres liés au protocole SMTP tels les listes noires ou de réputation, cadences de connexion, respect du protocole SMTP, ...
- **filtres de contenu** - ce sont les filtres qui ne tiennent compte que du contenu de la partie DATA du dialogue SMTP, c'est-à-dire, tout ce que l'on peut déduire des entêtes et du corps du message.

L'évaluation des filtres protocolaires est souvent plus délicate puisqu'elle dépend d'événements qui ne sont présents que dans des conditions réelles de fonctionnement ou alors de paramètres qui ne sont pas toujours observables tels le classement des messages refusés par une liste de réputation d'adresses : dans ce cas précis, on ne peut pas évaluer le taux de faux positifs. L'évaluation de ces filtres dépend souvent de la mise en place un dispositif spécifique au type de filtrage. Dans le cas d'une liste de réputation, par exemple, cela consiste à pouvoir recevoir tous les messages, même ceux qui seraient refusés à tort ou à raison, en y ajoutant une information de marquage pertinente. À partir du moment où l'évaluation de ces paramètres est possible, les principes sont les mêmes que pour l'évaluation des filtres de contenu. Certains aspects de l'évaluation de ces filtres sont cités dans cet article, mais nous nous attachons surtout à l'évaluation des filtres de contenu.

2 Évaluation de filtres anti-spam – principes et challenges

Le processus de filtrage de spam est généralement représenté par le modèle de la Figure 1. Une suite de messages est présentée, de façon séquentielle, à un classificateur qui les traite, aussi de façon séquentielle, et les attribue à une classe (ham ou spam). Le destinataire peut garder ou détruire le message, selon le classement donné par le filtre. Il peut aussi, vérifier si le classement du filtre est correcte et retourner des informations d'erreur (ou pas) au filtre, de façon à mettre à jour les modèles sur lesquels le filtre base son jugement (c'est *l'apprentissage en ligne*). Il convient de considérer que le retour

²Précision globale – ce paramètre est connu, dans la littérature en anglais par le terme « Accuracy ». A ne pas confondre avec « Précision ».

³A noter que la frontière entre ces deux types de filtrage est parfois floue.

d'information de la part du destinataire n'est pas systématique, n'est pas immédiate et, en plus, peut être erronée.

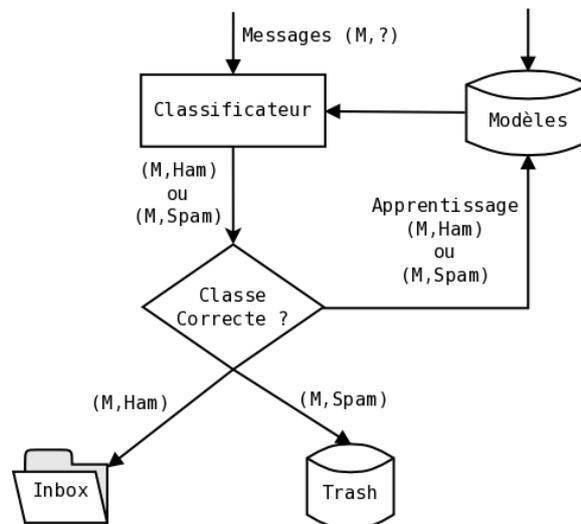


Figure 1: Modèle générique d'un processus de filtrage de spam

2.1 Les exigences des méthodes d'évaluation

Le protocole d'évaluation d'une « boîte noire », que ce soit un filtre anti-spam ou pas, doit satisfaire quelques contraintes :

- Lorsque l'évaluation est faite hors ligne, elle doit reproduire aussi fidèlement que possible l'environnement et l'utilisation réelle du filtre – cette contrainte évidente, n'est pas souvent suivie. Une erreur souvent commise est de considérer que le spam est un processus qui n'évolue pas dans le temps. Dans un autre exemple, un filtre de contenu placé après un filtre de *greylisting* n'est pas soumis au même type de contenu, ni en quantité, ni en qualité, que s'il n'y avait pas de filtrage par *greylisting*.
- Le test doit évaluer des critères d'efficacité significatifs et pertinents – parmi la multitude de critères d'efficacité utilisés dans les différentes applications de classement (diagnostique médical, classement de textes, recherche documentaire, ...) tous ne sont pas significatifs pour le classement du spam.
- Le test doit être effectué dans un environnement contrôlé – Il faut que l'environnement d'évaluation soit maîtrisé et stable. Ceci ne veut pas dire que l'on ne puisse pas tenir compte des événements inattendus, à condition qu'ils soient quantifiables en valeur, datés et dont l'impact sur la mesure soit aussi quantifiable et limité.
- Les résultats de l'évaluation doivent être statistiquement valides – il faut non seulement pouvoir évaluer des valeurs mais aussi la confiance que l'on peut accorder à cette valeur, par exemple, une estimation de l'erreur de la mesure.
- Le test doit être reproductible – comme toute expérimentation scientifique, tout doit être mis en œuvre pour que les résultats soient vérifiables. Cela implique que toutes les conditions de test doivent être datées enregistrées. Et sauf si l'évaluation n'a qu'un intérêt interne, les ensembles de données utilisés pendant l'expérimentation doivent être accessibles, ou alors assimilables (de façon vérifiable) à des données publiques.

2.2 Les aspects à prendre en compte lors de l'évaluation

L'évaluation d'un filtre anti-spam doit prendre en compte un nombre important d'aspects particuliers à ce type de dispositif. Dans cette section nous en énumérons quelques uns [10].

2.2.1 Non stationnarité

Ceci est un aspect particulièrement important du filtrage de spam, et pas souvent pris en compte dans les procédures d'évaluation.

Le trafic de la messagerie électronique varie selon trois aspects :

- La répartition des classes - La figure 2 montre la variation du taux de spam à l'entrée de l'École des Mines de Paris. À noter qu'il s'agit d'une mesure après filtrage par le *greylisting* (la dynamique serait plus importante si la mesure avait été faite avant). On peut nettement distinguer l'activité de nuit et de weekend. Cet exemple montre une variation à court terme, mais on remarque, sur une période plus longue (une année, par exemple), une variation plus ou moins périodique dans les hams, en apparence liée aux vacances et fêtes et une variation non périodique des

spams, plutôt liée à des événements autres (début ou fin d'activité d'un spammeur, débranchement de McColo, événement lié à une célébrité, ...).

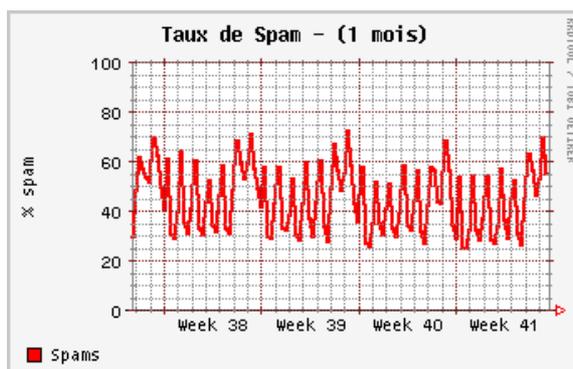


Figure 2: Taux de Spam à l'entrée de l'Ecole des Mines de Paris, après filtrage par greylisting

- La répartition à l'intérieur de chaque classe – ceci est plutôt vrai pour les spams. Rien ne garanti que la répartition par genre (pornographie, arnaques, médicaments, ...) soit constante dans le temps.
- Une variation qualitative à l'intérieur des messages – À l'intérieur des hams, les gens changent peu souvent leur façon d'écrire : les variations sont lentes. À l'intérieur des spams, au contraire, pour déjouer les filtres, les spammeurs changent souvent le contenu et la présentation des messages.

Ces trois aspects, pris individuellement, font que le spam n'est pas un processus statique (stationnaire) : il évolue dans le temps. Les deux derniers aspects font que le respect de l'ordre chronologique des messages est importante. Dans une évaluation hors ligne la soumission de messages doit se faire dans l'ordre chronologique.

NOTE : Dans la bibliographie concernant l'évaluation de filtres anti-spam, il y a une certaine confusion au sujet des expressions « en ligne » et « hors ligne », puisque on dit que, vu qu'il s'agit d'un processus non stationnaire, le filtrage de spam est un « processus en ligne », en opposition à un processus hors ligne, où l'ordre des messages n'a pas d'importance et l'apprentissage est entièrement effectué avant tout classement [11]. Ceci prête confusion avec une «évaluation en-ligne», faite sur un système en production, en opposition à une «évaluation hors-ligne», lorsque l'évaluation est faite par simulation de l'environnement réel. Sauf mention contraire, ce papier concerne « l'évaluation hors ligne d'un processus en ligne (non stationnaire) ».

Des méthodes courantes telles la validation croisée [12], plus adaptée aux processus stationnaires, ne peuvent donc pas être utilisés. Dans cette méthode, l'ensemble d'exemples est divisé en K sous ensembles et, utilisant à tour de rôle, K-1 pour l'apprentissage et le restant pour l'évaluation.

2.2.2 Évaluation en ligne et aspects temps réel

Lorsqu'il s'agit d'un système en production, il est souvent souhaitable de pouvoir faire de l'évaluation en ligne du filtre anti-spam. Mais, malgré le réalisme des conditions d'évaluation, certaines difficultés sont incontournables et relèvent, pour la plupart, du manque de maîtrise de l'environnement - les conditions de fonctionnement sont imposées par l'environnement avec ses aléas et l'évaluation de tous les paramètres de fonctionnement peut ne pas être possible. Cela fait qu'une évaluation n'est valable qu'au moment où elle est effectuée et difficilement reproductible. Ainsi, ces évaluations sont plus intéressantes pour évaluer la stabilité, dans le temps, d'un filtre utilisé en production, que pour évaluer son efficacité effective.

Des événements externes imprévus peuvent modifier l'environnement et impacter l'efficacité de filtrage. Il faut être capable de détecter ces changements de façon à pouvoir expliquer les résultats, avec enregistrement de toutes les informations pertinentes. Parfois, ces enregistrements peuvent permettre de rejouer, hors ligne, l'évaluation.

Certains changements dans l'environnement restent souhaitables et normaux (par exemple, la mise à jour d'une liste noire), mais rendent la mesure difficilement reproductible. Par contre, un événement tel l'inaccessibilité d'un serveur DNS d'une liste noire peut ne pas être détectable et fausser la mesure.

Il est aussi possible que certaines caractéristiques des filtres ne puissent pas être évaluées. Un exemple est le taux de faux positifs d'une liste, parce que que des messages sont rejetés sans qu'ils puissent être jugés de la classe à attribuer. Assez souvent, cette information n'est disponible que si les destinataires se rendent compte que certains messages ont été perdus, ce qui n'arrive pas systématiquement.

Si le but est l'évaluation intrinsèque de filtrage, une alternative souvent acceptable est d'enregistrer tous les événements survenus pendant une période donnée, accepter tous les messages (même ceux qui seraient refusés par une liste noire) et puis les dérouler dans un environnement contrôlé.

Un challenge d'évaluation en ligne de filtres anti-spam a été proposé lors des éditions 2007 et 2008 de CEAS [13] avec des résultats intéressants, mais pas encore satisfaisants. Le challenge de filtrage en ligne devra reprendre en 2010.

Certains organismes font l'évaluation en ligne de filtres anti-spam [14] [15]. Ce sont des évaluations intéressantes, mais vu qu'elles sont faites avec du trafic réel, elles ne sont valables qu'au moment où elles ont été réalisées : on ne peut donc pas comparer deux filtres qui n'ont pas été évalués en même temps.

2.2.3 Interactions avec le destinataire

Le retour d'information sur le classement (correct ou pas) est l'interaction la plus intéressante à étudier, puisqu'elle concerne directement les modes d'apprentissage du classificateur. Les résultats de filtrage influent et modifient le comportement de l'utilisateur et vice versa. La connaissance de cette interaction est un sujet encore ouvert et de ce fait, on privilégie les évaluations hors ligne pour étudier ces interactions. Des exemples de scénarios étudiés sont :

- Retour immédiat – c'est le cas le plus simple de l'utilisateur idéal qui retourne l'information d'exactitude de classement du message immédiatement après que le filtre l'a classé. Cette information est aussitôt intégrée à l'apprentissage du filtre, avant même l'arrivée du message suivant.
- Retour différé – c'est le cas de l'utilisateur qui ne lit son courrier que occasionnellement ou à des intervalles fixes et donc le retour d'information se fait quelque temps après le filtrage, et provoque une baisse d'efficacité (à évaluer).
- Retour sélectif – c'est le cas de l'utilisateur qui ne retourne que l'information d'exactitude de classement que pour une partie des messages (par exemple, que les spams non détectés ou que quand il a du temps libre).
- Uniquement les erreurs – l'utilisateur ne retourne que les classements en erreur. Éventuellement, il peut s'agir que des erreurs sur une seule classe. Les utilisateurs regardent plus souvent les dossiers de messages légitimes que celui de spams et, de ce fait, remarquent plutôt les spams non détectés que les messages légitimes filtrés à tort.
- Retour sur demande – (aussi appelé Apprentissage Actif), il s'agit de la possibilité pour le classificateur de demander au destinataire le classement correct des messages dont le score se trouve dans une zone d'indécision, information qui retournera avec un certain retard.
- Erreurs de l'information de retour – le destinataire n'étant pas un classificateur parfait, les informations de classement qu'il retourne peuvent ne pas être exactes. Assez souvent, ceci est modélisé comme du bruit qui contribue à baisser l'efficacité des filtres.

2.2.4 Interactions avec l'expéditeur

Les interactions du type *greylisting*, challenge/réponse, *captchas*, sont assez difficiles à simuler et à évaluer. Dans le cas du *greylisting*, par exemple, on ne peut pas évaluer dans un intervalle donné, la proportion de messages qui (avec le bon classement) ont été refusés grâce au *greylisting*, et le délai réel associé ne peut pas être évalué pour tous les messages acceptés. Ce qui se fait habituellement est comparer les résultats d'évaluation faites avec et sans la fonctionnalité en question. mais la validité de la démarche est relative puisque, le trafic de messagerie n'étant pas un processus stationnaire, on ne peut pas garantir que les deux évaluations ont été faites dans les mêmes conditions.

2.2.5 Corpus de messages utilisé pour l'évaluation

Des données différentes impliquent des résultats différents : rien ne peut assurer qu'un même filtre aura la même l'efficacité pour deux destinataires différents, ou pour le même destinataire à des instants différents. Donc, il est important que des comparaisons soient faites ayant comme référence le même ensemble de messages utilisé pour l'apprentissage et pour le test. S'il s'agit de comparer ses résultats avec des résultats obtenus par quelqu'un d'autre, il est essentiel que le corpus de messages soit le même.

La constitution d'un corpus public de messages n'est pas une tâche simple : difficile à convaincre quelqu'un de publier la totalité de ses messages. Les premières publications de résultats de recherche ont été effectués avec des corpus de messages en provenance de listes de diffusion ouvertes (Ling-Spam [16], PU [17] et SpamAssassin [18]) ou alors des messages dont les attributs ont été déjà extraits et chiffrés. Il n'y avait aucune assurance que ces messages étaient représentatifs d'une boîte aux lettres réelle.

Après la faillite de la société Enron en 2001, l'utilisation, pour des besoins de recherche, des messages trouvés dans les serveurs de mail de l'entreprise a été autorisé. C'est ainsi que le corpus de messages de TREC a été constitué [18] [19]. Néanmoins, malgré l'intérêt de ce corpus, les résultats que l'on peut obtenir ne sont pas généralisables. Martins et Cormack [20] [21] ont montré qu'il ne pouvait pas être utilisé pour étudier les caractéristiques temporelles des filtres anti-spam. Ce

corpus a aussi d'autres faiblesses : il a été constitué à partir de messages qui ont presque 10 ans et qui sont, pour la plupart, en texte pur (pas de partie HTML), les messages légitimes sont tous en langue anglaise, ne permettant donc pas d'étudier le filtrage dans un environnement multi-langue ou dont la langue principale n'est pas l'anglais.

3 TREC - Spam Track

TREC (« Text REtrieval Conference ») [22] est un workshop sponsorisé par deux départements de l'administration américaine : le NIST (National Institut of Standards and Technologie) et le DoD (Department of Defense). Le but étant de promouvoir, annuellement, des échanges entre chercheurs de tous azimuts, autour de thèmes novateurs en rapport avec la problématique de la recherche documentaire.

Le thème Spam a été traité pendant trois ans, de 2005 à 2007. Le but était simple : évaluer si les filtres anti-spam étaient vraiment efficaces. À l'époque, les produits commerciaux étaient proposés avec une « sauce secrète » sans aucune possibilité de réelle vérification de leurs caractéristiques (et c'est toujours le cas). Les filtres anti-spam libres évoquaient une excellente efficacité, mais avec des mesures empiriques on non répétables [5] [7]. La communauté de la recherche en intelligence artificielle étudiait le filtrage de spam utilisant un modèle abstrait d'apprentissage supervisé avec des vecteurs d'attributs déjà extraits des messages, comme si l'algorithme de classement était la partie la plus critique du filtrage anti-spam.

Le principe de TREC a toujours été d'améliorer l'état de l'art grâce à des évaluations, mesures et comparaisons effectuées conjointement à l'aide de jeux de données publiquement disponibles. La problématique du spam est différente dans le sens où il s'agit d'un problème de classement en ligne, sur un flot de données non stationnaire et pratiquement tous les jeux de données sont privés.

Cormack et Lynam [23] ont développé une boîte à outils avec une interface permettant de soumettre au filtre, pour classement, une suite de messages selon un ordre chronologique précis et d'enregistrer les réponses. Les indicateurs d'efficacité pertinents ainsi que leur signification statistique sont évalués à partir des réponses du filtre. Deux corpus de messages ont été utilisés : un public constitué à partir des messages récupérés après la faillite de la société Enron et un autre privé constitué à partir des messages reçus par un utilisateur sur une période de 8 mois.

L'efficacité des filtres participants à TREC Spam Track s'est amélioré à chaque année, par rapport à l'année précédente. En 2005, Bratko a démontré que les modèles à compression étaient plus performants que les filtres pseudo-bayésiens de l'époque. En 2006, le filtre proposé par Assis, utilisant des digrammes orthogonaux épars pour modéliser les messages, avec un seuil d'apprentissage, a été le plus performant. En 2007, les filtres basés sur des classificateurs à régression logistique et SVMs (Support Vector Machines) et utilisant des n-grams au lieu de mots pour modéliser les messages ont été les plus performants [24].

4 Métriques d'évaluation

Dans cette section, nous présentons quelques uns des indicateurs les plus pertinents pour l'évaluation de filtres anti-spam de contenu.

4.1 Tables de contingence et valeurs associées

La table de contingence (ou table de co-occurrence) est un outil souvent utilisé lorsqu'on souhaite étudier les relations entre deux variables pouvant prendre des valeurs discrètes (ou des catégories). Dans notre cas, les variables sont, dans les colonnes, le classement réel (aussi connu par « *Gold Standard* ») et dans les lignes, le résultat du filtre. La somme de chaque colonne donne le nombre réel d'éléments dans chaque classe et celle de chaque ligne donne le nombre d'éléments vus par le classificateur dans chaque classe. Les différents rapports que l'on peut extraire de la table permettent de définir des critères d'efficacité, plus ou moins pertinents selon le type d'application.

Un exemple numérique de table de contingence est présenté dans la Table 1, avec les variables associées : VN (vrais négatifs : le nombre de hams vus par le filtre comme étant des hams), FN (faux négatifs : spams vus comme des hams), VP (vrais positifs : spams vus comme des spams) et FP (faux positifs : hams vus comme des spams).

	Vrai Ham	Vrai Spam
Classement Ham	9982 (VN)	334 (FN)
Classement Spam	18 (FP)	39666 (VP)

Table 1 - Table de Contingence – classement réel versus classement du filtre

- **Probabilité à priori des classes** – il ne s'agit pas d'un indicateur de l'efficacité de filtrage, mais d'un paramètre de contrôle indiquant les conditions de fonctionnement du filtre.

$$P_{spam} = (VP + FN) / (VP + FN + VN + FP)$$

$$P_{ham} = (VN + FP) / (VP + FN + VN + FP)$$

- **Taux d'erreur par classe** (faux positifs et faux négatifs) – il s'agit de la fraction du nombre d'objets d'une catégorie classés par erreur dans l'autre classe.

$$FPR = FP / (VN + FP)$$

$$FNR = FN / (VP + FN)$$

Ce sont des critères pertinents et aussi plus intuitifs dans les applications de classement de spams. Ils ont l'avantage de ne pas dépendre des probabilités à priori de chaque classe (rapport ham/spam).

A noter que, pour évaluer ces taux d'erreur, certains fournisseurs utilisent la quantité totale de messages et non pas la quantité par classe : cela permet de présenter des meilleurs résultats, parfois un ordre de grandeur plus bas.

- **Taux de bon classement** (vrai positifs et vrai négatifs ou sensibilité et spécificité)

$$VPR = VP / (VP + FN) = 1 - FNR$$

$$VNR = VN / (VN + FP) = 1 - FPR$$

- **Précision et Rappel** (Precision and Recall) – la précision indique la proportion de spams parmi les messages détectés comme étant du spam, tandis que le rappel est le ratio entre le nombre de spams détectés comme tel et le nombre total de spams.

$$Precision = VP / (VP + FP)$$

$$Recall = VP / (VP + FN)$$

Ces deux critères ont leur origine dans les applications de recherche documentaire et on les trouve parfois dans les résultats d'évaluation de filtres anti-spam, mais ils ne sont pas pertinents pour cette application, en particulier la précision, à cause de leur dépendance des probabilités à priori des classes.

- **Précision Globale** (Accuracy) – il s'agit du taux total d'erreurs, les deux classes confondues.

$$Accuracy = (VP + VN) / (VP + FN + VN + FP)$$

Ce indicateur est souvent mentionné, mais il n'a d'intérêt que quand les classes sont plus ou moins symétriques [Provost98], ce qui n'est pas le cas du spam où les probabilités à priori des classes, les taux d'erreurs usuels d'opération ainsi que le coût des erreurs sont souvent très différents.

- **Taux d'erreur pondérés** – Compte tenu de l'asymétrie des classes, plusieurs méthodes de calcul de taux d'erreur pondérés par un « coefficient de risque », associé à chaque classe ont été proposées. Actuellement le consensus est que ces indicateurs ne présentent pas d'intérêt : d'une part, même si l'on accepte l'asymétrie des classes, il n'y a pas de sens de choisir une valeur plutôt qu'une autre et d'autre part, à l'intérieur d'une même classe, le risque associé à une erreur de classement n'est pas uniforme. Il est préférable, par exemple, d'afficher les valeurs pertinents et de laisser à chacun l'interprétation.

4.2 R.O.C. (Receiver Operating Characteristic) et 1-ROCA

Dans la section précédente nous avons vu quelques indicateurs déductibles des tables de contingence. Ces indicateurs ont l'inconvénient d'être spécifiques à un point d'opération particulier du classificateur – une valeur de seuil - et ne disent rien sur l'efficacité du filtre à d'autres points d'opération.

Certains filtres présentent le résultat de façon binaire - ham/spam ("*hard classifiers*") tandis que d'autres sous la forme d'une valeur numérique de score ("*soft classifiers*"). Lorsque cette valeur est disponible, on peut définir des seuils tels que les messages dont le score est inférieur seront classés dans une catégorie et ceux dont le score est supérieur dans l'autre. Si l'on trace les courbes de taux d'erreurs de classement en fonction du score, on obtient une courbe similaire à celle de la Figure 3.

Le diagramme ROC (*Receiver Operating Characteristic*) [25] est un outil que l'on trouve dans des domaines tels le diagnostic médical et la radiologie, mais qui a ses origines dans les problèmes de détection radar (d'où son nom). Il s'agit de la courbe paramétrique du *taux de vrai positifs* (1 - taux de faux négatifs) en fonction du *taux de faux positifs* paramétrée par la valeur de seuil.

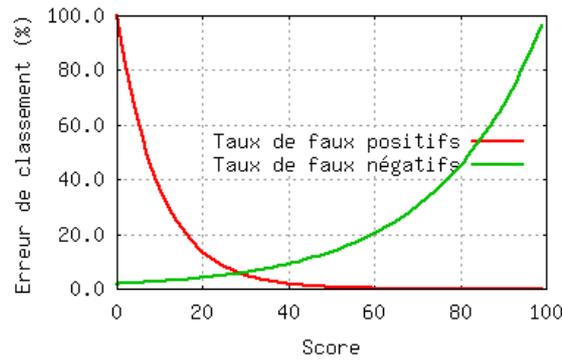


Figure 3: Erreur de classement en fonction du score

L'exemple dans la Figure 4 présente un morceau (coin en haut à gauche) des courbes ROC de deux configurations différentes (unités de segmentation du texte) du même filtre. Pour tracer ces courbes il suffit de soumettre un ensemble de messages au classificateur et noter pour chaque valeur de score les taux de faux positifs et de faux négatifs.

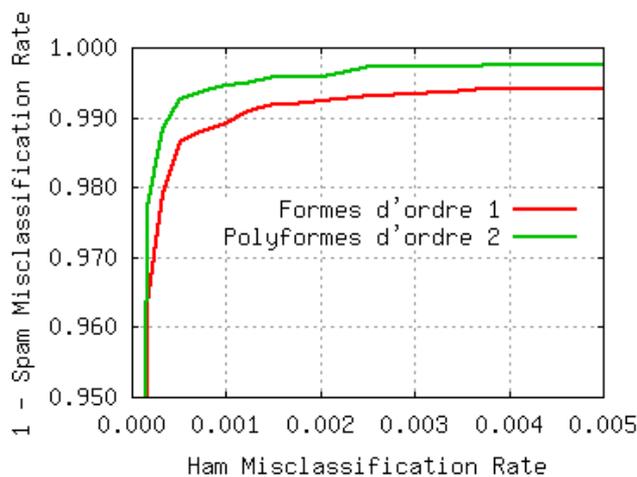


Figure 4: Exemple de courbes ROC pour le même filtre, utilisant des unités de segmentation différentes

Les courbes ROC ont des caractéristiques visuelles que l'on peut interpréter facilement :

- la courbe ROC est entièrement comprise dans le carré de sommets opposés (0,0) et (1,1) puisque les variables que l'on représente dans les deux axes sont des probabilités,
- un filtre idéal (au cas où il existerait un) ne commet pas d'erreurs et donc sa courbe ROC se confond avec les côtés gauche et supérieur de ce carré,
- le segment reliant les sommets (0,0) et (1,1) définit une zone d'incertitude où l'efficacité du filtre est identique à un choix aléatoire. Si la courbe se trouve au dessous ou à droite de ce segment, il vaut mieux faire un choix aléatoire plutôt qu'utiliser le filtre, ou alors prendre le choix inverse.
- pour comparer deux points d'opération, il suffit de sélectionner celui qui est plus vers le haut et vers la gauche.
- pour comparer deux classificateurs, le meilleur est celui dont la courbe est la plus vers le haut et vers la gauche – dans l'exemple de la figure 4, il est préférable d'utiliser des polyformes d'ordre 2 que des formes simples.

La courbe ROC a aussi une autre propriété importante qui caractérise l'efficacité intrinsèque du classificateurs. Si l'on prends, au hasard, deux objets à classer, un dans chaque catégorie, la probabilité que leurs scores soient ordonnés dans le bon ordre de classement est égale à la surface au dessous de la courbe. En général, c'est la valeur complémentaire à 1 qui est utilisée, et que l'on trouve sous les sigles (1-AUC) ou (1-ROCA). Dans l'état actuel, la valeur de ce indicateur pour les filtres anti-spam les plus performantes se situe entre 0,01 et 0,05 %.

Cette courbe donne aussi un critère permettant de comparer un classificateur binaire ("*hard classifier*") et un autre avec un classificateur non binaire ("*soft classifier*"). Il suffit de placer le couple (VPR, FPR) dans le même graphique que celui avec qui on souhaite comparer. Si ce point se trouve à l'intérieur de la courbe, cela veut dire qu'il est probablement moins bon que l'autre.

4.3 L.A.M. (Logistic Average Misclassification – Erreur Moyen Logistique)

La courbe ROC est intéressante puisqu'elle permet d'évaluer globalement un classificateur indépendamment du point d'opération. Les inconvénients de ROC sont le besoin d'avoir accès à la valeur du score attribué à chaque message filtré et le manque de relation directe entre la valeur de (1-AUC) et les taux d'erreur.

Lors de TREC 2005, il a été remarqué [10] que le LAM (Erreur moyen logistique) varie peu avec la valeur du score.

Le LAM est défini par :

$$LAM = \text{logit}^{-1}\left(\frac{\text{logit}(FPR) + \text{logit}(FNR)}{2}\right)$$

où *logit* est la fonction erreur logistique (et son inverse) :

$$\text{logit}(x) = \log\left(\frac{x}{1-x}\right) \quad \text{et} \quad \text{logit}^{-1}(y) = \frac{1}{1+e^{-y}}$$

Cet indicateur a une valeur qui est toujours intermédiaire entre le taux de faux positifs et le taux de faux négatifs et est proche de la moyenne géométrique, pour des petites valeurs.

En réalité, la propriété de stabilité de la valeur de LAM n'est pas le seul intérêt de cet indicateur. Les taux d'erreur et de bon classement se situent dans un intervalle fermé [0,1], avec des valeurs souvent proches des valeurs extrêmes (et donc écrasées). L'utilisation de la fonction *logit* permet de transposer cet intervalle fermé en un intervalle ouvert $[-\infty, +\infty]$, et l'échelle logarithmique évite l'écrasement lorsque les deux valeurs à moyenner diffèrent d'un ordre de grandeur ou plus. Par exemple, si le taux de faux négatifs est de 2% et celui de faux positifs est de 0,1 %, une amélioration d'un facteur 10 dans ce dernier n'est pas ressentie si l'on utilise une moyenne arithmétique, alors qu'il apparaît clairement avec l'utilisation de la moyenne logistique.

4.4 Sensibilité au bruit

Pour les filtres à apprentissage, le bruit est un facteur qui peut détériorer la capacité de généralisation. Des exemples de bruit sont les retours d'information erronées par les destinataires (bruit dans le classement), des spams incluant des "mots savants" (bruit dans les attributs) ou alors, ce que l'on remarque assez souvent actuellement, la reprise d'un message généré par une *newsletter* avec remplacement des URLs et des images les plus visibles par celles du site du spammeur.

La modélisation du bruit n'est pas aussi simple, puisque les erreurs ne sont pas distribuées de façon homogène sur les différents types de messages et d'attributs. Néanmoins, la simulation d'un bruit de classement homogène permet déjà d'avoir une idée de la résistance du filtre au bruit. Pour cela, il suffit de modifier, aléatoirement, les classements des messages soumis à l'apprentissage, avec des taux entre, disons, 1 % et 10 %.

4.5 Signification statistique de la mesure

L'objectif de la mesure d'efficacité d'un filtre dans une expérimentation particulière est de prédire son efficacité dans des situations similaires. La confiance que l'on peut accorder à cette mesure dépend de deux facteurs : le degré de réalisme avec lequel l'expérimentation représente les situations à prédire et les erreurs dues aux aléas de la mesure.

Ces erreurs sont données, en général, par un intervalle de confiance à 95 %, c'est-à-dire, si on répète N fois la mesure, l'intervalle à l'intérieur duquel le résultat doit se trouver 95 % du nombre de mesure. Vu que les jugements du filtre permettent d'avoir la distribution empirique des scores, une méthode adaptée à ce type d'évaluation est le "*Bootstrap*" (c'est la méthode utilisée dans TREC). Pour une description détaillée de cette méthode, voir le livre de Efron et Tibshirani [26].

Outre l'évaluation de l'intervalle de confiance, il est important, autant que possible, de s'assurer que la quantité de messages à tester est suffisante pour la plage de valeurs à mesurer. Par exemple, tester un filtre sur seulement 1000 messages n'est pas suffisant si le taux de faux positifs attendu est de l'ordre de 0,1 %.

5 Évaluation Hors ligne – La boîte à outils d'évaluation de TREC Spam Track

Dans les sections précédentes, nous avons décrit les contraintes à respecter lors de l'évaluation d'un filtre anti-spam ainsi que quelques paramètres intéressants à évaluer. Idéalement, si l'on veut pouvoir comparer deux filtres différents ou deux configurations différentes du même filtre, il faut pouvoir soumettre les filtres en étude aux mêmes conditions d'évaluation, si possible de façon automatisée.

Les évaluations faites dans TREC ont utilisé une boîte à outils développée par Cormack et Lynam [23] avec ces caractéristiques. Cette boîte à outils était constituée d'une partie générique de contrôle, et d'une interface spécifique à chaque filtre permettant à la partie contrôle d'envoyer des commandes au filtre et à recevoir les réponses.

Cette interface implémente quatre commandes :

- **initialize** - initialise le filtre et démarre tous les services nécessaires à son fonctionnement le rendant prêt à classer des nouveaux messages ou à intégrer des messages au modèle d'apprentissage
- **classify emailfile** - classe un message (fichier emailfile) et retourne le résultat : la classe (ham/spam) et le score (un nombre réel).
- **train class emailfile** - ajoute le message au modèle d'apprentissage et retourne son classement avant apprentissage.
- **finalize** - termine le filtre, ainsi que tous les services démarrés lors de la phase d'initialisation. Après avoir envoyé cette commande au filtre.

La partie contrôle lit un fichier avec la spécification de l'évaluation, lance les actions d'évaluation dans un ordre précis, enregistre les résultats et, à la fin, déduit les indicateurs recherchés décrits dans ce papier, à partir des résultats du filtre en étude.

La partie spécifique au filtre peut être facilement programmée dans les cas où le filtre possède une interface en ligne de commande ou des mécanismes standard de communication sous UNIX (pipes, ...). Dans les cas des filtres où le mode de communication est au travers d'un protocole du type SMTP, on peut encore trouver une solution, grâce à un outil de soumission et d'un outil de réception de messages.

Cette boîte à outils, distribuée sous licence GPL, peut être utilisée sans modification ou alors être personnalisée selon les besoins de chacun.

6 Discussion et conclusions

L'évaluation des filtres anti-spam est devenu un besoin aussi bien pour ceux qui développent des filtres, mais aussi pour ceux qui les utilisent.

Dans cet article nous avons passé en revue les points importants à prendre en compte lors de l'évaluation d'un filtre anti-spam, ainsi que les indices d'efficacité pertinents les plus courants. Ces indices d'efficacité ne sont d'aucune valeur s'ils ne sont pas accompagnés des informations permettant de savoir comment ils ont été évalués. Certains outils de base tels la table de contingence nous permettent d'évaluer l'efficacité à des points d'opération précis, tandis que d'autres indices tels la *1-ROCA* et *LAM* sont des critères d'agrégation, permettant d'avoir une idée plus globale du fonctionnement du filtre. Néanmoins, l'efficacité d'un filtre est une information à plusieurs dimensions et on ne peut pas toujours se contenter de l'information sur un tout petit nombre d'indices.

L'ensemble qui a été décrit dans cet article est largement inspiré de la méthodologie d'évaluation de filtres anti-spam utilisé dans les évaluations qui ont eu lieu pendant les trois années de réalisation de TREC Spam Track. L'auteur tient à remercier Gordon Cormack, de l'Université de Waterloo par les nombreux échanges concernant le filtrage de spam et, en particulier, l'évaluation des filtres.

Bibliographie

- [1] M. Sahami, S. Dumais, D. Heckerman, and E. Horvitz. *A bayesian approach to filter junk e-mail*. In AAI-98 Workshop in Learning for Text Categorization , 1998.
- [2] P. Graham. *A plan for spam*. <http://www.paulgraham.com/spam.html>, 2002.
- [3] Vaderetro – *Documentation commerciale* : http://www.vade-retro.com/doc/plaquette_vaderetro.pdf – accessible le 16 octobre 2009
- [4] Ironport – *Documentation (Datasheet)* http://www.ironport.com/pdf/ironport_anti-spam_datasheet.pdf – accessible le 16 octobre 2009

- [5] W. S. Yerazunis, *The Spam-Filtering Accuracy Plateau at 99.9% Accuracy and How to Get Past It*, MIT Spam Conference, January 2004 ([TR2004-091](#))
- [6] G. V. Cormack and T. R. Lynam, *TREC 2005 Spam Track Overview* – <http://plg.uwaterloo.ca/~gvcormac/trecspamtrack05>, 2005.
- [7] G. V. Cormack and T. R. Lynam, *Online supervised spam filter evaluation*. *ACM Transactions on Information Systems* 25, 3 (Jul. 2007), 11.
- [8] J. A. Zdziarski, *Ending Spam*, No Starch Press, 2005,
- [9] Provost, F. J., Fawcett, T., and Kohavi, R. 1998. *The Case against Accuracy Estimation for Comparing Induction Algorithms*. In *Proceedings of the Fifteenth international Conference on Machine Learning* (July 24 - 27, 1998). J. W. Shavlik, Ed. Morgan Kaufmann Publishers, San Francisco, CA, 445-453.
- [10] Gordon V. Cormack. *Email Spam Filtering : A systematic review*, volume 1. Now Publishers, 2008
- [11] N. Cristianini and J. Shawe-Taylor, *Support Vector Machines and other kernel based learning methods*, Cambridge University Press, 2000
- [12] Tom Mitchell, *Machine Learning*, McGraw-Hill, 1997
- [13] *The CEAS 2008 live spam challenge*. <http://www.ceas.cc/2008/challenge/challenge.html>, 2008
- [14] OpusOne Inc, - *Spam Testing Methodology* - <http://www.opus1.com/www/whitepapers/spamtestmethodology.pdf> - March 2007
- [15] Virus-Bulletin - *VBSpam Testing Methodology* - <http://www.virusbnt.com/vbspam/methodology/index> - April 2009
- [16] Ling-Spam Corpus, http://www.aueb.gr/users/ion/data/lingspam_public.tar.gz
- [17] PU Corpus, <http://www.aueb.gr/users/ion/data/PU123ACorpora.tar.gz>
- [18] SpamAssassin Corpus, <http://spamassassin.apache.org/publiccorpus/>
- [19] G. V. Cormack and T. R. Lynam, *Spam corpus creation for TREC*. in CEAS 2005 : The Second Conference on E-mail and Anti-spam, 2005
- [20] G. V. Cormack and J. M. Martins da Cruz. *On the relative age of spam and ham training samples for email filtering*. In SIGIR '09 : Proceedings of the 32nd international ACM SIGIR Conference on Research and Development in Information Retrieval, pages 744–745, New York, NY, USA, 2009
- [21] J. M. Martins da Cruz and G. V. Cormack. *Using old spam and ham samples to train email filters*. In Proc. CEAS 2009 – Sixth Conference on Email and Anti-Spam, Mountain View, CA, 2009.
- [22] TREC, <http://trec.nist.gov>
- [23] G. V. Cormack and T. R. Lynam. *TREC spam filter evaluation toolkit*. <http://plg.uwaterloo.ca/~gvcormac/jig/>.
- [24] G. V. Cormack, *Communication personelle*, 2009
- [25] T. Fawcett. *An Introduction to ROC Analysis*. *Pattern Recogn. Lett.*, 27(8) :861–874, 2006.
- [26] B. Efron and R. J. Tibshirani. *An Introduction to the Bootstrap* . Chapman and Hall, New York, 1994