

Status report del gruppo di lavoro GARR “sec-mail”

Roberto Cecchini¹, Fulvia Costa², Alberto D’Ambrosio³, Giacomo Fazio⁴, Antonio Forte⁵, Matteo Genghini⁶, Michele Michelotto², Ombretta Pinazza⁷, Alessandro Spanu⁵, Alfonso Sparano⁸

Abstract

Verrà presentato un report dell’attività del gruppo di lavoro GARR “sec-mail”.

I punti principali possono essere sintetizzati in: studio di metodologie per migliorare l’efficienza dei metodi di rivelazione dello spam (in particolare SpamAssassin); stesura di un documento di “best practice” sulla posta elettronica e sistemi di autenticazione del dominio del mittente dei messaggi di posta elettronica, in particolare SPF.

Introduzione

Oggetto della presentazione è un report dell’attività del gruppo di lavoro GARR “sec-mail”, nato su proposta di Roberto Cecchini durante il workshop GARR di Roma del Novembre 2003.

Il gruppo si occupa dei problemi collegati alla sicurezza informatica della posta elettronica nella rete GARR. In particolare sono stati trattati i seguenti argomenti:

- spam;
- controllo di virus e worm diffusi attraverso la posta elettronica;
- definizione di “best practices” per la configurazione dei servizi di posta e per la configurazione dei filtri di protezione dei mail server;
- autenticazione del server del mittente (parzialmente legato al problema degli spam).

Da fine 2004 il gruppo dispone di un’area web nella quale vengono resi disponibili i risultati delle sperimentazioni e la documentazione prodotta .

Il gruppo si è concentrato soprattutto nel tuning di SpamAssassin per migliorare l’identificazione degli spam mediante l’uso di filtri bayesiani, regole aggiuntive non standard e metodi basati su identificazione distribuita di mass e-mailing. Il gruppo si è dotato di alcuni server DCC che ha reso disponibili come servizio pilota alla rete GARR.

Controllo dello spam

La base comune: SpamAssassin

L’attività del gruppo si è concentrata nel capire come fosse possibile migliorare l’efficienza di SpamAssassin sia nel ridurre i falsi negativi, ma soprattutto evitando la comparsa di falsi positivi.

SpamAssassin ha un insieme di regole euristiche rispetto alle quali tutti i mail devono essere confrontati. Se un mail in arrivo ha un riscontro positivo con una certa regola, a questo mail viene sommato il valore della regola (che può anche essere negativo). Se il mail supera una certa soglia, di default fissata a 5, il mail vie-

¹ INFN, Firenze

² INFN, Padova

³ INFN, Torino

⁴ INAF, Palermo

⁵ INFN, Roma

⁶ IASF, Bologna

⁷ INFN, Bologna

⁸ Università di Salerno

ne considerato come spam. A questo punto l'amministratore può decidere di cancellare il mail o di metterlo in una cartella dell'utente separata dalla normale Inbox, per consentirgli di decidere come comportarsi. Una scelta molto comune è quella di cambiare il soggetto per aiutare l'utente a filtrare il mail per i mail sopra soglia.

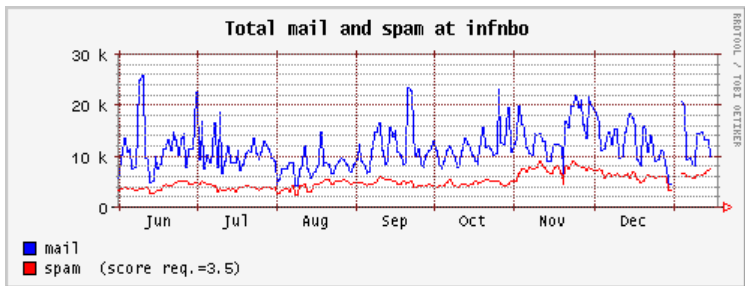
Un metodo banale per aumentare il numero di mail identificati come spam è quello di abbassare la soglia oppure di aumentare il valore di punteggio di alcune regole. Questo metodo è molto pericoloso perché abbassare la soglia vuol dire aumentare la probabilità di avere falsi positivi. Anche alzare il punteggio di alcuni test è sconsigliabile perché il punteggio di ogni test viene assegnato in base al punteggio di tutti gli altri test.

Abbiamo quindi cercato di implementare metodi indipendenti per aumentare l'efficacia di separazione tra la curva di distribuzione degli ham e quella degli spam.

Un metodo molto efficace è quello del filtro *bayesiano*. Con questo metodo un certo insieme di email identificati da un umano come sicuri e tipici spam e ham vengono catalogati dal filtro. Il filtro ricava due tabelle delle parole (in realtà dei token perché l'analisi viene fatta anche nell'header dei mail) più frequenti negli spam e negli ham. Dopo questa fase di apprendimento, ad ogni mail in arrivo viene assegnato un punteggio di probabilità di spam in base alla prevalenza dei token ham o spam. Questo metodo è in gran parte indipendente da quello basato sulle regole perché lavora anche sui mail buoni e inoltre è personalizzato sul target di ham di ogni mailservet.

Abbiamo riscontrato alcune debolezze in questo metodo, in particolare i database delle parole ham e spam con il tempo invecchiano e quindi può essere necessario aggiornarli. Alcuni raccomandano di implementare le opzioni di autolearning ma la nostra esperienza ci ha fatto capire che è molto facile che proprio con l'autolearning il sistema peggiori più rapidamente a causa dei mail che vengono catalogati male. In particolare gli spammer di tanto in tanto inviano mail con insiemi di parole casuali progettati con lo scopo di ingannare i filtri bayesiani. Il rimedio consiste nel correggere il database bayesiano (centralizzato o per utente) in base ai feedback degli utenti (riclassificazione di falsi positivi e falsi negativi).

Stiamo realizzando un sistema di monitoraggio dell'efficienza dei vari test impiegati da SpamAssassin (cfr. la figura e la tabella seguente, riferite ad una sola sede, ma in corso di estensione ad altre).



Valori numerici riferiti all'ultimo log: Sun Jan 16 23:59:42 2005

Test	score	hit	hit%	eff	eff%
WS_URI_RBL	1.500	5163	1.6	62	0.98
JP_URI_RBL	2.500	4902	7.5	67	1.06
OB_URI_RBL	2.200	4806	6.0	56	0.89
HTML_MESSAGE	0.001	3910	1.8	1	0.02
SPAMCOP_URI_RBL	4.000	3762	9.5	90	1.42
PYZOR_CHECK	0.322	3583	6.7	1	0.02
DCC_CHECK	1.806	3457	4.7	135	2.13

MIME_HTML_ONLY	0.100	2563	0.5	0	0.00
AB_URI_RBL	3.000	1870	9.6	22	0.35
RAZOR2_CHECK	0.899	1399	2.1	11	0.17
RAZOR2_CF_RANGE_51_100	1.552	1390	2.0	23	0.36
HTML_IMAGE_ONLY_02	0.100	1094	7.3	0	0.00
HTML_40_50	0.474	769	2.2	4	0.06
HTML_90_100	1.073	753	1.9	0	0.00
FROM_ENDS_IN_NUMS	0.869	643	0.2	1	0.02
HTML_FONT_BIG	0.100	578	9.1	0	0.00
HTML_50_60	0.183	534	8.4	0	0.00
MIME_HTML_ONLY_MULTI	1.101	530	8.4	2	0.03
SAVE_UP_TO	0.100	484	7.7	0	0.00
HTML_IMAGE_ONLY_06	0.100	483	7.6	0	0.00
PRIORITY_NO_NAME	0.100	423	6.7	0	0.00
ONLINE_PHARMACY	0.100	402	6.4	0	0.00
FORGED_YAHOO_RCVD	3.900	381	6.0	49	0.77

Dal qualche mese è disponibile la versione 3.0 di SpamAssassin. In un primo momento le prestazioni in termini di rilevazione degli spam sembravano decisamente inferiori a quelle delle versioni 2.6x. Le ultime prove con la versione 3.02, invece, mostrano un netto miglioramento: test su di un campione di circa 1000 mail, analizzati sia con la versione 2 sia con la 3, hanno mostrato una riduzione dei falsi negativi dal 20 all'1% e dei falsi positivi dallo 0.8 allo 0.2%.

RBL

Un altro sistema di tagging indipendente è quello basato sulle RealTime Block List. Il mittente del mail viene confrontato con una query dns con un database di noti spammer. Questi database sono mantenuti da gruppi di utenti che raccolgono le segnalazioni di Internet Service Provider che usano ospitare spammer, o che permettono open relay, Proxy male configurati, invio di posta da indirizzi dinamici. Le block list hanno spesso dato origine a polemiche nel passato perché alcune erano molto lente nel togliere i siti che si mettevano in regola oppure non erano accurate nel controllare che il sito denunciato fosse veramente in difetto. Per questo motivo è bene usare i punteggi delle RBL come punteggio aggiuntivo alle altre regole e non basarsi esclusivamente sulla presenza in una RBL per identificare il mail come Spam.

Abbiamo trovato particolarmente utile una RBL particolare che si basa sulla raccolta dei siti che sono puntati dalle URL all'interno dei mail di tipo spam. Questo perché mentre spesso gli header dei mail non permettono di capire da dove veramente arriva lo spam, all'interno del contenuto utile del messaggio ci devono essere link utili agli spammer.

Razor, Pyzor e DCC

Questi sistemi si basano su server distribuiti che contano i mail ricevuti per determinare la probabilità che siano di tipo UBE.

Razor

Razor [RAZ] si basa sulla sottomissioni di mail identificati da umani come spam. Un meccanismo di punteggi basato sulla correttezza delle segnalazioni di spam o di revoca di uno spam garantisce che gli spammer non possano alterare il meccanismo di segnalazione. Razor si basa però su un protocollo chiuso e non è facile entrare in queste rete collaborativa. In questo momento l'interrogazione dei server Razor è gratuita.

Pyzor

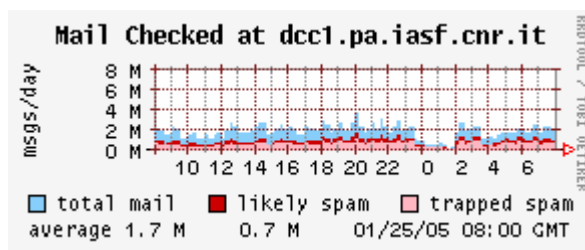
Pyzor [PYZ] è un tentativo di replicare la rete Razor con software open source e con partecipazione aperta.

DCC

DCC invece si basa su un meccanismo leggermente diverso. I server DCC contano automaticamente i mail bulk, cercano di togliere gli elementi variabili e tenendo solo gli elementi fissi di un mail, generano un checksum. I server DCC poi si scambiano con un meccanismo di flooding queste segnalazioni. Un mail server che vuole essere cliente di un server DCC può richiedere per ogni mail in arrivo la probabilità che il mail con quel checksum sia di tipo UBE. Un Server DCC risponde sia a client registrati che a client anonimi. Viene data priorità ai client registrati

DCC è un sistema aperto in cui viene incoraggiata la partecipazione. Tutti i siti dei partecipanti al working group hanno cominciato ad usare i client DCC e in tre sedi sono stati installati dei server. Il primo server DCC in Italia è stato installato presso l’INAF di Palermo e poi in seguito a questa esperienza sono stati installati presso l’INFN di Torino e di Roma. Il nostro obiettivo è di aumentare ancora il numero dei server DCC per accogliere le richieste di siti GARR che ne facciano richiesta. La speranza è di migliorare l’efficienza del metodo sia diminuendo i tempi di risposta, sia perché il database si arricchirebbe di dati sui mail di spam “nostrani”, attualmente non presenti nei server in funzione. Non c’è nessuna decisione ufficiale di offrire un servizio che non sia l’attuale su base *best effort*.

Verranno presentati dati sull’uso dei server e sull’efficienza del metodo di rivelazione. Ad esempio, il server DCC di Palermo, nella giornata del 24/1/04, ha ricevuto circa 350k richieste di checksum e 1.2M report dagli oltre 7000 clienti autorizzati.



DSPam

DSPam [DSP] è un sistema di rilevazione spam che si propone come un’alternativa a SpamAssassin. Si basa solo su tecniche statistiche molto sofisticate. Gli autori si propongono un’efficienza migliore del 99.9%, anche se nei nostri test preliminari non solo non siamo riusciti a raggiungere questa percentuale, ma anzi abbiamo ottenuto risultati inferiori a quelli di SpamAssassin. C’è da dire che in dspam è molto importante che la qualità del training sia elevata, dal momento che il sistema si basa esclusivamente su metodi statistici.

Stiamo valutando la possibilità di utilizzarlo come uno dei test di SpamAssassin.

Best Practice

Una parte dell’attività è stata dedicata alle cosiddette “best practice”, cioè ai consigli per migliorare la sicurezza dei servizi di mail.

Tra i punti più importanti identificati.

- Consentire dai router di frontiera o dal firewall l’accesso alla porta 25 incoming solo verso il mail server ufficiale del dominio per evitare che macchine della LAN facciano da relay.
- Limitare l’accesso alla porta 25 in uscita ai soli mailer ufficiali, per evitare che eventuali virus o worm sulla rete locale possano spedire messaggi di posta (al momento un comportamento tipico).

Le porte 587 e 465 vanno lasciate aperte per permettere agli utenti roaming l'utilizzo del proprio MTA.

- Permettere l'utilizzo dall'esterno, previa autenticazione, del proprio MTA: essenziale se si vogliono applicare metodi di controllo del mittente (ad es. SPF, vedi in seguito).
- Configurare l'eventuale prodotto Antivirus nel mail server in modo che non mandi la segnalazione al finto mittente di un mail infetto dal momento che il mittente viene quasi sempre falsificato (spoofed).
- Attivare il meccanismo di **GreetPause** (da sendmail 8.13), grazie al quale una connessione SMTP viene rifiutata dal server se chi chiede la connessione non aspetta la risposta di greeting “220”. Questo è un tipico comportamento dei software di spamming e dei virus. Quasi sempre il software dello spammer, dopo aver ricevuto il rifiuto della connessione, non riprova dopo un certo timeout, come fanno gli MTA regolari.

Uno dei membri del gruppo (D'Ambrosio) è coautore di un manuale per l'installazione e la configurazione di un servizio di posta elettronica con **sendmail [MAN]**, che può essere considerato un buon punto di partenza.

Autenticazione del server del mittente

L'autenticazione del server del mittente è un altro argomento importante, anche perché parzialmente collegato al problema spam (molti mail di questo tipo hanno infatti il mittente falsificato) e legato a questo la necessità di autenticare gli utenti esterni alla LAN per permettere l'accesso ai mail server autorizzati. degli utenti roaming che devono poter accedere ai servizi di posta elettronica dall'esterno della propria LAN.

Sono state proposte diverse metodologie, senza che si sia riusciti ad unificarle in un RFC. Tra le varie proposte si sono distinte due tecnologie, la prima nota come SPF **[SPF]**, la seconda come Sender-ID **[SID]**. Entrambe non servono direttamente per la lotta contro lo spam ma servono solo ad autenticare il server del mittente di una e-mail. Il concetto di base utilizzato da questi sistemi è che un utente di un dato dominio può mandare mail solo tramite il server autorizzato del proprio dominio. Questo impedisce l'invio di mail con mittente falsificato, ma soprattutto evita che eventuali computer infettati spediscono mail contenenti spam o virus.

Abbiamo deciso di verificare gli eventuali benefici dell'uso di SPF. Un sito che vuole diventare SPF compliant deve pubblicare nei propri record DNS i nomi delle macchine autorizzate ad inviare posta a suo nome. Il massimo dei risultati si avrebbe qualora la grande maggioranza dei siti decidesse di pubblicare la lista dei server autorizzati, cosa cui siamo ancora ben lontani, tuttavia, dal momento che alcuni grandi ISP pubblicano già il record SPF, è possibile sfruttare questa informazione per modificare il punteggio di SpamAssassin. È bene ricordare che l'ISP che decide di pubblicare il record SPF deve informare i propri utenti roaming che non possono inviare e-mail “firmandosi” con il nome del dominio di appartenenza qualora dovessero utilizzare indirizzi IP esterni al proprio dominio, perché fallirebbe il controllo SPF. Quindi, prima di pubblicare il record SPF, è necessario istruire i propri utenti affinché utilizzino sempre i server autorizzati per inviare mail. Per tale motivo si deve permettere il relay ai propri utenti, utilizzando gli opportuni meccanismi di autorizzazione (ad es. via password o certificato).

Durante la verifica, effettuata in un dominio reale, si è constatato che su 650K mail, ricevute nell'arco di mese, il 12% circa provenivano da mittenti i cui amministratori avevano già pubblicato informazioni SPF per i propri domini. A questo valore si può aggiungere un altro 20% di mail che rappresenta la quantità di posta inviata da computer appartenenti al dominio preso in esame (e quindi autorizzati secondo le specifiche SPF), arrivando quindi ad un ragguardevole 32% di mail dotate di informazioni SPF.

Bibliografia

[DCC]

<http://www.rhyolite.com/anti-spam/dcc/>

[DSP]

<http://www.nuclearelephant.com/projects/dspam/>

[MAN]

Giorgio Bar, Alberto D’Ambrosio, Franca De Giovanni, *Manuale di installazione di un servizio di posta elettronica completo di filtri anti-virus e anti-spam* (Versione Pre-2)

[PYZ]

<http://pyzor.sourceforge.net/>

[RBL]

<http://www.webopedia.com/TERM/R/RBL.html>

[RAZ]

<http://razor.sourceforge.net/>

[SID]

<http://www.senderid.org/>

<http://www.microsoft.com/mscorp/twc/privacy/spam/senderid/default.aspx>

[SPF]

<http://spf.pobox.com/>

[WSM]

<http://www.garr.it/WG/sec-mail/>

Glossario

Spam	Mail che riceviamo ma che non siamo interessati a ricevere.
Ham	I mail buoni, usato in contrapposizione a Spam
UCE	Unsolicited Commercial Email. Mail pubblicitari non richiesti e quindi si presume non graditi. Indipendentemente dal fatto che il mittente mandi un solo mail ad un solo destinatario o ne mandi invece a migliaia
UBE	Unsolicited Bulk Email. Mail mandati a migliaia di utenti. Possono anche non essere commerciali ma essere per esempio di tipo politico, o per sostenere una causa, o anche solo per testare la reale esistenza dei destinatari.
Falsi Negativi	Mail di tipo Spam che i sistemi automatici antispam non riescono a identificare come Spam
Falsi Positivi	Mail di tipo Ham che i sistemi automatici antispam identificano erroneamente come Spam.

Biografie degli autori

Roberto Cecchini.

Laureato in Fisica, è responsabile del Servizio Calcolo e Reti della Sezione INFN di Firenze, dal 1999 è responsabile del servizio di sicurezza informatica della rete GARR (GARR-CERT), dal 1998 gestisce la Certification Authority dell’INFN (INFN CA).

Fulvia Costa.

Lan Manager e APM presso la sezione INFN di Padova.

Alberto D’Ambrosio

Perito Elettronico Industriale, System Administrator c/o Servizio Calcolo INFN (dal 2000 della Sezione INFN di Torino, dal 1992 al 2000 dei Laboratori Nazionali del Gran Sasso), in precedenza anali-

sta/programmatore nel campo dell'automazione industriale presso aziende private.

Giacomo Fazio

Informatico, System and Network Administrator presso lo IASF sezione di Palermo. Responsabile per la Posta Elettronica dell'intero IASF CNR. Responsabile dei Servizi di Calcolo IASF Palermo. Dal 1991 presso l'IFCAI, poi divenuto IASF ora INAF. In precedenza programmatore per i gruppi di Ricerca operanti nel campo dell'Astrofisica.

Antonio Forte

Antonio Forte, perito industriale informatico, system administrator presso i servizi di calcolo delle sezioni INFN di Roma2 (1996-1998), Torino (1998-2003) e Roma (dal 1/1/2004).

Matteo Genghini

Laureato in Ingegneria Elettronica, dal 2002 è responsabile dei servizi informatici del Centro di Calcolo dell'Istituto IASF di Bologna. In precedenza programmatore in campo multimediale presso aziende private.

Michele Michelotto

Fisico, ha lavorato presso il CERN di Ginevra e l'INFN di Legnaro e Padova principalmente nell'elaborazione offline in esperimenti di fisica delle alte energie. Dal 1997 Tecnologo all'INFN di Padova come responsabile del locale servizio calcolo.

Ombretta Pinazza

Laureata in Fisica nel 1993, ha svolto attività di ricerca presso la Facoltà di Ingegneria dell'Università di Bologna ed è stata poi ricercatrice al Tesre-CNR di Bologna; dal 1999 è tecnologo informatico presso la Sezione INFN di Bologna.

Alessandro Spanu

Alessandro Spanu, system e network manager presso la Sezione INFN di Roma dal 1982. Dal 2000 è responsabile del Servizio Impianti Calcolo e Reti della Sezione INFN di Roma.

Alfonso Sparano

Laureato in Ingegneria Informatica, è amministratore di servizi informatici offerti dal Centro Elaborazione Dati al personale e ai Dipartimenti dell'Università degli Studi di Salerno.