

Indice

1	Introduzione	1
1.1	Ranking dei motori di ricerca	3
1.1.1	Metodi di ranking endogeno	4
1.1.2	Metodi di ranking esogeno	5
1.2	Web spam	6
1.2.1	Tecniche di boost	7
1.2.2	Term Spamming	7
1.2.3	Link Spamming	9
1.2.4	Tecniche di hiding	11
1.2.5	Click Spamming	12
2	Survey sulle tecniche di spam detection	14
2.1	Tecniche basate sul contenuto	14
2.1.1	Prime feature per identificare lo spam	14
2.1.2	Utilizzo di un classificatore per combinare le feature	23
2.1.3	Language model per rilevare lo spam	24
2.1.4	Spam detection sulla base degli argomenti di una pagina web	29
2.1.5	Altre tecniche	33
2.2	Tecniche basate sul grafo	34
2.2.1	Metodi classici per identificare lo spam web usando il grafo	34
2.2.2	Metodi per identificare spam farm	39

Elenco delle figure

1.1	Tassonomia delle tecniche boost	8
1.2	Tipi di pagine nel web per uno spammer	10
1.3	Esempio di una spamfarm	11
1.4	Tecniche di hiding	12
2.1	Occorrenze dello spam classificate per dominio all'interno del dataset descritto in [17]	17
2.2	Occorrenze dello spam classificate per lingua all'interno del dataset descritto in [17]	17
2.3	Prevalenza di spam sulla base del numero di parole per pagina . . .	18
2.4	Prevalenza di spam sulla base del numero di parole all'interno dei titoli delle pagine	18
2.5	Prevalenza di spam sulla base della frazione di parole di una pagina che sono tra le 200 più frequenti parole nel corpus	19
2.6	Prevalenza di spam sulla base della lunghezza media delle parole per pagina	20
2.7	Prevalenza di spam sulla base della quantità di testo delle ancore delle pagine	21
2.8	Prevalenza di spam sulla base della frazione di contenuto visibile . .	21
2.9	Prevalenza di spam sulla base del rapporto di compressione	22
2.10	Esempio di classificatore	23

2.11 Istogramma della divergenza KL tra testo delle ancore e il contenuto della pagina puntata basato sul dataset WEBSPPAM-UK2006 utilizzato in [16]	25
2.12 Istogramma della divergenza KL tra testo intorno alle ancore e il contenuto della pagina puntata basato sul dataset WEBSPPAM-UK2006 utilizzato in [16]	26
2.13 Istogramma della divergenza KL tra termini degli URL e il contenuto della pagina puntata basato sul dataset WEBSPPAM-UK2007 utilizzato in [16]	26
2.14 Istogramma della divergenza KL tra testo delle ancore e titolo della pagina puntata basato sul dataset WEBSPPAM-UK2007 utilizzato in [16]	27
2.15 Istogramma della divergenza KL tra testo intorno alle ancore e titolo della pagina puntata basato sul dataset WEBSPPAM-UK2006 utilizzato in [16]	27
2.16 Istogramma della divergenza KL tra termini nell'URL e titolo della pagina puntata basato sul dataset WEBSPPAM-UK2006 utilizzato in [16]	28
2.17 Istogramma della divergenza KL tra titolo e contenuto della pagina basato sul dataset WEBSPPAM-UK2006 utilizzato in [16]	28
2.18 Istogramma della divergenza KL tra testo delle ancore e i meta tag della pagina basato sul dataset WEBSPPAM-UK2006 utilizzato in [16]	29
2.19 Distribuzione degli argomenti pesati per pagine spam e normali . . .	30
2.20 Prevalenza di spam relativa alla misura della diversità degli argomenti basata sulla varianza	31
2.21 Prevalenza di spam relativa alla misura di diversità degli argomenti basata sulla semantica	32
2.22 Prevalenza di spam relativa alla misura di diversità sulla massima semantica	33
2.23 Algoritmo di trustrank	36

Capitolo 1

Introduzione

Questa tesi ha come obbiettivo lo studio e l'analisi delle tecniche di spam detection attualmente esistenti ed in particolare delle tecniche online. Nella prima parte le tecniche verranno classificate sulla base dei segnali che utilizzano. Successivamente verranno eseguiti dei test per valutare alcuni algoritmi di spam detection offline eseguiti durante la fase di crawling. Ed infine verranno presentati e discussi i risultati ottenuti. Attualmente sono poche le tecniche online di spam detection, ovvero tecniche che rilevano lo spam durante la fase di crawling. Infatti quasi tutti i metodi tentano di fare il crawling dell'intera porzione di web di interesse e successivamente classificare le pagine in classi (dove di norma le classi sono due: *spam* oppure *non spam*).

Il fenomeno del web spam è sempre più presente all'interno del web: questo è dovuto al fatto che gli utenti tendono ad esaminare solo i primi risultati calcolati dai motori di ricerca e quindi se un sito compare tra i primi n risultati può avere un ritorno economico maggiore, legato alla quantità di traffico che viene generata per quel sito. Da uno studio del 2005 [10] si stima che la perdita finanziaria mondiale causata dallo spam è di circa 50 miliardi di dollari mentre nel 2009 tale perdita è salita a 130 miliardi di dollari, come descritto in [11]. Risulta ovvio, quindi, perché recentemente tutte le più grandi compagnie di motori di ricerca hanno identificato nel recupero di informazioni non pertinenti una delle priorità da risolvere.

Le conseguenze del web spam possono essere riassunte come segue [20]:

- la qualità delle ricerche è compromessa penalizzando i siti web legittimi;
- un utente potrebbe perdere la fiducia sulla qualità di un motore di ricerca e passare con facilità all'utilizzo di un altro;
- inoltre i siti spam possono essere usati come mezzo per malware, pubblicazione di contenuto per adulti e attacchi di tipo “fishing”. Una prova tangibile si può vedere in [4], dove gli autori hanno eseguito l'algoritmo di *PageRank* su 100 milioni di pagine e hanno notato che 11 sui primi 20 risultati erano composti da siti con contenuto per adulti.

Queste considerazioni evidenziano che nella progettazione di un motore di ricerca occorre tenere conto delle pagine che potrebbero portare al mal funzionamento del motore stesso.

Il lavoro prodotto da questa tesi sarà utilizzato per essere integrato all'interno di un web crawler distribuito ad alte prestazioni per il futuro sviluppo di un modulo di spam detection. L'esigenza di tale modulo è sorta a seguito dello sviluppo, presso il Dipartimento, di un crawler chiamato *BUBiNG*, altamente configurabile ma privo al momento di qualunque forma di rilevazione di siti e contenuti malevoli. Il problema è estremamente interessante sia dal punto di vista teorico che da quello pratico: infatti, sebbene siano numerose le tecniche descritte in letteratura per la determinazione di spam (usando come segnali sia il contenuto che la struttura dei link), è sorprendentemente scarso l'insieme di tali tecniche che possono essere usate on-line, cioè durante il crawl. Il problema diventa ancora più complesso se si aggiungono considerazioni legate ai vincoli di spazio di memoria disponibile e al tempo di calcolo.

In letteratura il processo di spam detection viene eseguito quasi sempre offline, subito dopo la fase di crawling. Tale processo è integrato nella fase di ranking dei motori di ricerca:

- crawling dell'intero web;
- fase di spam detection;
- indicizzazione.

Il motivo per cui la fase di spam detection viene eseguita dopo la fase di crawling è perché si utilizza il grafo risultante per determinare le pagine spam.

Partendo da queste considerazioni, in questa tesi verranno effettuate delle analisi per determinare se il processo di spam detection possa essere eseguito durante la fase di crawling ovvero al momento in cui il crawler esegue il “fetch” di una pagina per determinare “on the fly” se la pagina è buona o ha un contenuto malevolo.

1.1 Ranking dei motori di ricerca

Prima di spiegare i vari metodi con cui si possono creare pagine web spam e i vari metodi utili ad identificarlo, è necessario capire come i motori di ricerca siano capaci di valutare la rilevanza di una pagina web per una determinata query.

In linea di massima un sistema di reperimento di informazioni, ovvero un motore di ricerca, è dato da una collezione documentale D (un insieme di documenti) di dimensione N , da un insieme Q di interrogazioni e da funzione di ranking ($r : Q \times D \rightarrow R$) che assegna ad ogni coppia formata da un’interrogazione e un documento un numero reale. L’idea è che a fronte di un’interrogazione a ogni documento venga assegnato un punteggio reale: i documenti con punteggio nullo non sono considerati rilevanti, mentre quelli a punteggio non nullo sono tanto più rilevanti quanto più il loro punteggio è alto. In particolare i metodi di ranking si dividono in *endogeni* ed *esogeni*. I primi metodi fanno uso del contenuto del documento per valutarne la rilevanza mentre i secondi fanno uso di una struttura esterna, ad esempio il grafo composto dai collegamenti ipertestuali tra le pagine web; questo non implica che i metodi esogeni non possano fare uso del contenuto della pagina (per esempio il testo delle ancore). I criteri si dividono ulteriormente in statici (o indipendenti dall’interrogazione) e dinamici (o dipendenti dall’interrogazione). Nel primo caso il punteggio assegnato a ciascun documento è fisso e indipendente da un’interrogazione q mentre nel secondo il punteggio assegnato a ciascun documento è dipendente da un’interrogazione q .

Tra i metodi endogeni sono di maggiore importanza *tf-idf* e *BM25* mentre tra quelli esogeni i più diffusi in letteratura sono *PageRank* e *HITS*.

1.1.1 Metodi di ranking endogeno

I metodi di ranking endogeno utilizzano il contenuto di una pagina per assegnarle un punteggio. Possono essere anch'essi statici o dinamici (cioè dipendere o meno da un'interrogazione). L'algoritmo usato dai motori di ricerca per fare il rank delle pagine web basandosi sui campi di testo usa varie forme del *tf-idf*. Il *tf-idf* è un metodo di ranking endogeno dinamico che utilizza il contenuto di una pagina per assegnarle un punteggio. Il *tf-idf* è composto da due misure più semplici: la *Term Frequency* e la *Inverse Document Frequency*. Il primo metodo assegna a un documento d il punteggio dato dalla somma dei conteggi dei termini t dell'interrogazione che compaiono nel documento stesso. In questo modo documenti in cui i termini dell'interrogazione compaiono più frequentemente avranno un punteggio più elevato. Utilizzare solo questo metodo non conviene in quanto è facilmente manipolabile. Inoltre non tiene conto del fatto che alcuni termini occorrono più frequentemente non perché rilevanti, ma perché altamente frequenti all'interno di *ogni* documento (ad esempio le congiunzioni). Il secondo metodo è definito come l'inverso del numero di documenti nella collezione che contengono il termine t [14]. Più precisamente:

$$idf_t = \log \frac{N}{df_t} \quad (1.1)$$

dove N è il numero totale di documenti nella collezione e df_t il numero di documenti nella collezione dove il termine t occorre. La combinazione del *tf* ed dell'*idf* produce una misura composta che permette di normalizzare il peso dei termini. Il *tf-idf* di un documento d rispetto a una query q è calcolato su tutti i termini t in comune come:

$$tf - idf(d, q) = \sum_{t \in d \text{ and } t \in q} tf(t, d) \cdot idf(t) \quad (1.2)$$

Con il *tf-idf* gli spammer possono avere due obiettivi: o creare pagine rilevanti per un gran numero di query o creare pagine molto rilevanti per una specifica query. Il primo obiettivo può essere ottenuto includendo un gran numero di termini distinti in un documento; il secondo, attraverso la ripetizione di determinati termini nel documento.

Un altro metodo di ranking endogeno è *BM25* [19], basato sul *modello probabilistico*.

1.1.2 Metodi di ranking esogeno

Uno dei metodi esogeni statici è *PageRank* descritto in [18]. *PageRank* usa le informazioni portate dai link in entrata (*inlink*) per determinare un punteggio globale di importanza di una pagina. Esso assume che esista un legame tra il numero di *inlink* di una pagina p e la popolarità della pagina p . L'importanza di una pagina web p , utilizzando *PageRank*, è legata al numero di pagine web che puntano ad essa e all'importanza di tali pagine web. Questo concetto è mutualmente rinforzante ovvero l'importanza di una certa pagina influenza ed è influenzata dall'importanza delle altre pagine [8].

PageRank è basato sulla passeggiata naturale del grafo del web G . Più precisamente, la passeggiata viene perturbata nel seguente modo: fissato un parametro α tra 0 e 1, a ogni passo con probabilità α si segue un arco uscente, e con probabilità $1 - \alpha$ si sceglie un qualunque altro nodo del grafo utilizzando una qualche distribuzione v , detta vettore di preferenza (per esempio, uniforme). Assumendo che non esistano pozzi, la matrice di transizione della catena è quindi rappresentata dalla combinazione lineare:

$$\alpha G + (1 - \alpha)1v^T \quad (1.3)$$

dove G è la matrice della passeggiata naturale su G . Il fattore α è detto fattore di attenuazione di norma è impostato a un valore di 0,85.

Un altro metodo esogeno usato per il ranking delle pagine è *HITS* (*Hyperlink Induced Topic Distillation*) introdotto in [12]. Differentemente da *PageRank* esso assegna due punteggi di importanza a ogni pagina: uno di *hubbiness* e uno di *autorevolezza*. L'intuizione dietro a *HITS* è che invece di un singolo punteggio di importanza esista un concetto di pagina *autorevole*, cioè pagina con contenuto pertinente e interessante, e di *hub*, cioè pagina contenente numerosi collegamenti a pagine autorevoli. I due concetti si rinforzano *mutuamente*: una pagina autorevole è pun-

tata da molte pagine centrali, e una buona pagina centrale punta a molte pagine autorevoli.

Questo approccio considera che nel web ci sono due tipi di pagine: quelle che contengono dei contenuti per un determinato argomento (*authoritative*) e quelle che contengono tanti link a delle pagine *authoritative* che sono chiamate pagine *hub*. Le pagine *hub* sono utili per scoprire le pagine *authoritative* [15].

L'algoritmo lavora su un sottografo del web ottenuto a partire da un'interrogazione. La selezione del sottografo può essere fatta in vari modi, un modo è quello di prendere un certo insieme di risultati ottenuto da un motore di base e generare un sottografo sulla base di una query e delle pagine che puntano a quelle ottenute dalla query. Per questo sottoinsieme di pagine otteniamo una matrice di adiacenza A . I punteggi di *hub* e *authority* per tutte le pagine del sottoinsieme possono essere formalizzate dalla seguente coppia di equazioni:

$$\begin{cases} \vec{a}_{t+1} = A^T \vec{h}_t \\ \vec{h}_{t+1} = A \vec{a}_{t+1} \end{cases} \quad (1.4)$$

Può essere dimostrato che la soluzione ottenuta applicando iterativamente il sistema 1.4 converge rispettivamente al principale autovettore di AA^T e $A^T A$ [15][20].

1.2 Web spam

Con il termine web spamming si fa riferimento a tutti i metodi che tentano di manipolare gli algoritmi di ranking dei motori di ricerca per aumentare il valore di alcune pagine rispetto ad altre [8]. Dato il numero esorbitante di pagine che vengono create e pubblicate sul web, gli utenti competono per far comparire le proprie pagine tra le prime dei risultati di una query. Il fenomeno dello spamming o spamindexing ricade sulla qualità delle ricerche causando diversi problemi: indicizzazione di pagine che non sono utili, aumento del costo delle operazioni di query, malware e reindirizzamento verso contenuto per adulti; inoltre questo spinge gli utenti ad utilizzare altri motori di ricerca [20].

L'obiettivo dei motori di ricerca è di ottenere ottimi risultati per identificare tutte le pagine web che sono rilevanti per una specifica query e presentarle secondo

l'importanza che esse hanno. Di norma la rilevanza viene misurata attraverso la similarità testuale tra la query e le pagine mentre l'importanza è definita come la popolarità globale della pagina e a volte è inferita dalla struttura dei link [8]. Ci sono due categorie di tecniche associate al web spam [8]:

- tecniche boost che cercano di far avere più importanza o rilevanza a delle pagine
- tecniche hiding che sono metodi per nascondere le tecniche di boost all'utente dal browser, anche se alcuni autori incorporano queste tecniche fra quelle di boost.

1.2.1 Tecniche di boost

Le tecniche di boosting si dividono in: *Term Spamming* e *Link Spamming*. Con l'avvento degli algoritmi di ranking basati sulla struttura del grafo il *Term Spamming* è stato trascurato. In figura 1.1 è rappresentata una possibile tassonomia delle tecniche boost [8].

1.2.2 Term Spamming

Nel valutare la rilevanza testuale i motori di ricerca considerano dove i termini di una query compaiono in una pagina. Il tipo di punto all'interno della pagina è chiamato *campo*. I più comuni campi di testo per una pagina p sono: il body della pagina, il titolo, i meta tag nell'header HTML e l'URL della pagina. Inoltre viene considerato anche come *campo*, il testo delle ancore (il tag a) associate all'URL che puntano alla pagina p dato che descrive molto bene il contenuto della pagina. I campi di testo di p sono utilizzati per determinare la rilevanza di p rispetto ad una query (alcune volte i campi vengono pesati sulla base della loro importanza) e perciò chi fa *term spamming* utilizza tecniche di pesatura dei contenuti dei campi di testo in modo tale da aumentare l'efficacia dello spam [8]. Le tecniche di spamming possono essere raggruppate in base ai *campi* di testo dove viene fatto spamming. In base a questo distinguiamo [8]:

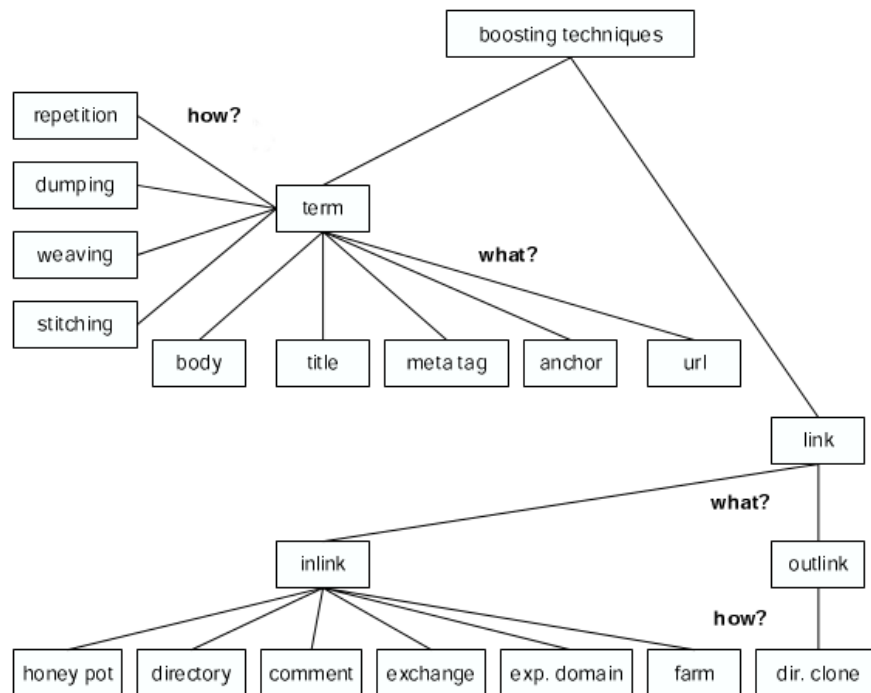


Figura 1.1: Tassonomia delle tecniche boost

- *Body Spam.* In questo caso lo spam è nel corpo del documento. Questo è lo spam più diffuso.
- *Title Spam.* Molti motori di ricerca danno molta importanza ai termini che compaiono nel titolo. Quindi ha senso includere termini di spam all'interno del titolo della pagina.
- *Meta Tag Spam.* I tag che compaiono nell'header sono molto frequentemente soggetti a spam. Per questo i motori di ricerca danno poca importanza a questi campi o non li considerano. Di seguito viene mostrato un esempio di questo tipo di spam.

```
<meta name="keyword" content="buy, cheap, cameras, lens,
    accessories, nikon, canon">
```

- *Anchor Text Spam.* I motori di ricerca assegnano un peso maggiore al testo nelle ancore perché pensano che esse contengano un riassunto del contenuto

della pagina. Perciò del testo di spam è incluso nel testo delle ancore dei collegamenti HTML di una pagina. In questo caso lo spamming non viene fatto sulla pagina cui si vuole far avere un rank più alto ma sulle pagine che puntano ad essa.

```
<a href="target.html">free, great deals, cheap,  
    inexpensive, cheap, free</a>
```

- *URL Spam*. Alcuni motori di ricerca dividono l'URL delle pagine in un insieme di termini che sono usati per determinare la rilevanza di una pagina. Per sfruttare questo metodo di ranking, gli spammer creano lunghi URL che includono una grande sequenza di termini spam, un esempio può essere: *buy-canon-rebel-20d-lens-case.camerasx.com*.

Queste tecniche possono essere utilizzate insieme o separatamente. Un altro modo per raggruppare queste tecniche si basa sul tipo di termini che vengono utilizzati nei campi di testo [8], possiamo avere:

- Ripetizione di uno o più specifici termini.
- Inclusione di molti termini generici per creare pagine rilevanti per molte query.
- Intreccio di vari termini all'interno della pagina.
- Creazione di frasi di senso compiuto per l'elaborazione di contenuti generati velocemente attraverso la concatenazione di frasi da fonti diverse.

1.2.3 Link Spamming

Il *link spamming* è un tipo di spam che fa uso della struttura dei link tra le pagine web per favorire il rank di una pagina target t . Come descritto in [8], per uno spammer ci sono tre tipi di pagine web: inaccessibili, accessibili (blog) e proprietarie (fig. 1.2). Le inaccessibili sono quelle che uno spammer non può modificare. Le accessibili sono pagine gestite da altri ma che possono essere modificate lievemente dallo spammer attraverso l'immissione di un post in un forum, in un blog o in portali di questo

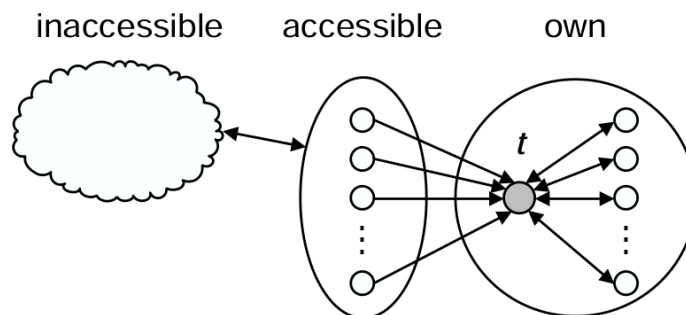


Figura 1.2: Tipi di pagine nel web per uno spammer

genere. Le proprietarie sono pagine su cui gli spammer hanno il pieno controllo. Il gruppo di pagine proprietarie è chiamato *spam farm*.

Molti motori di ricerca utilizzano due algoritmi per aumentare l'importanza basandosi sulle informazioni dei link: PageRank e HITS; sulla base di questi due tipi di algoritmi vengono definite due categorie principali di *link spamming*: *outgoing link spam* e *incoming link spam*. L'*outgoing link* è uno dei metodi più facili da implementare in quanto basta aggiungere alla propria pagina dei link ad altre pagine, che sono considerate buone, al fine di aumentare il punteggio di *hub*. Per la ricerca di link da includere nella pagina per cui si vuole incrementare il punteggio di *hub* si possono utilizzare delle directory che contengono liste di siti come DMOZ o Yahoo!. Queste directory organizzano i contenuti web in contenuti e in liste di siti relativi. Per quanto riguarda *incoming link*, ci sono diverse strategie che si possono adottare in modo tale da avere un numero elevato di link in entrata [8]:

- *Honeypot*: consiste nella creazione di un insieme di pagine aventi un contenuto interessante (per esempio una documentazione Linux) ma che contengono link nascosti alla pagina o alle pagine di cui si deve aumentare il valore di rilevanza.
- *Infiltrarsi in una directory web*: tecnica con cui i webmaster inseriscono nelle directory web link ai loro siti.
- *Postare link nei blog, forum e wiki*: consiste nell'includere URL a pagine di spam come parte di un commento.

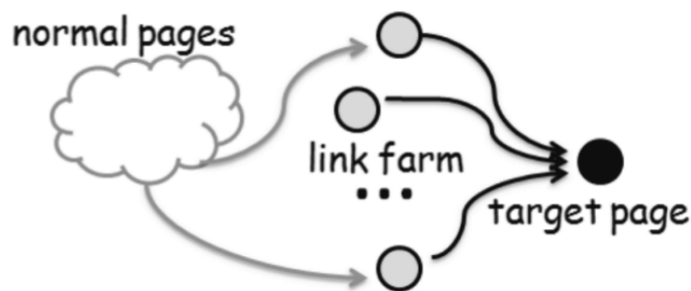


Figura 1.3: Esempio di una spamfarm

- *Scambio di link*: tecnica che consiste nello scambiare link con altre pagine di spam. Rappresenta una pratica comune tra chi fa spam ed infatti esistono blog completamente finalizzati all'incontro di spammer per lo scambio dei link.
- *Comprare domini scaduti*: tecnica che consiste nell'acquisto di domini scaduti e nella manipolazione delle pagine che puntano ancora ad esso al fine di aumentare il rank di una pagina target.
- *Creare una spam farm*: consiste nella creazione di un insieme di pagine che puntano ad una pagina spam, detta *target page*, e che hanno come obiettivo l'aumento della rilevanza di quest'ultima. In questo caso il valore di page rank aggregato delle pagine è propagato alla pagina target. Molte volte tale tecnica si integra con quella *honeypot*. Una delle forme più aggressive di integrazione honeypot - spamfarm è l'*hijacking* [20], dove gli spammer prima attaccano un sito con una buona reputabilità e poi usano questo come parte della loro link farm.

1.2.4 Tecniche di hiding

Le tecniche di hiding si possono classificare in: *content hiding*, *cloaking*, *redirection* (fig. 1.4) [8]. Nel *Content hiding* i termini o i link di spam vengono nascosti quando il browser visualizza una pagina, ad esempio utilizzando per i termini lo stesso colore dello sfondo e per i link non inserendo il testo all'interno delle ancore che indirizzano ad una pagina. Un'altra tecnica è quella di utilizzare degli script per

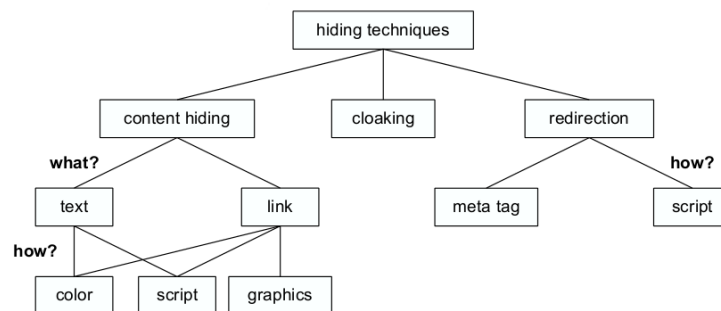


Figura 1.4: Tecniche di hiding

nascondere il contenuto. Il *Cloaking* sfrutta la possibilità di identificare se la richiesta di una pagina è fatta da un crawler o da un browser: dato un URL, il server spam restituisce un documento HTML diverso a seconda che la richiesta venga fatta da un crawler o da un browser. Quindi vengono distribuiti due contenuti diversi in base alla provenienza della richiesta al server spam. La rilevazione di un crawler può essere effettuata in due modi: o si mantiene in memoria una lista di indirizzi di crawler oppure si usa l'header della richiesta HTTP, controllando il campo user-agent e verificando che sia diverso dai più comuni browser (si ipotizza quindi che sia un crawler). Nell'esempio sotto, lo user-agent della richiesta HTML indica l'uso del web browser Chrome.

```

Mozilla/5.0 (X11; Linux x86_64)
AppleWebKit/537.36 (KHTML, like Gecko)
Chrome/32.0.1700.102 Safari/537.36

```

La *Redirection* è un'altra tecnica che reindirizza il browser ad un altro URL appena la pagina è caricata. Un esempio di redirection server-side è mostrato di seguito.

```

header("Location: http://www.example.com/");

```

1.2.5 Click Spamming

Un ultimo metodo per fare web spam è il *Click Spamming* [20]. I motori di ricerca utilizzano dati sul flusso dei click per regolare le funzioni di ranking, quindi gli

spammer generano clic fraudolenti per manipolare il comportamento di queste funzioni in modo tale da fare avere un rank migliore ai loro siti. Il metodo prevede che vengano fatte delle query e si clicchi sulla pagina di cui si vuole aumentare il rank. Tale metodo viene eseguito in modo automatico attraverso script che girano su diverse macchine per non fare sospettare il motore di ricerca delle numerose richieste provenienti da un'unica macchina [20].

Capitolo 2

Survey sulle tecniche di spam detection

Il capitolo illustra alcune tecniche di spam detection presenti in letteratura. Le tecniche verranno suddivise sulla base del tipo di segnali che vengono utilizzati: il contenuto, il grafo ed altri segnali (come ad esempio l'header delle richieste HTTP). Nella prima parte del capitolo verranno illustrate le tecniche basate sul contenuto, nella seconda parte le tecniche che fanno uso di grafi ed infine le tecniche che fanno uso di altri tipi di segnali.

2.1 Tecniche basate sul contenuto

2.1.1 Prime feature per identificare lo spam

Un metodo per identificare lo spam basandosi sul contenuto di una pagina web è quello di analizzare alcune proprietà (feature) delle pagine spam e confrontarle con le medesime proprietà di quelle non spam al fine di ottenere dei valori con cui stimare la natura della pagina web (spam o non spam). Alcune di queste proprietà, come descritto in [5], sono le seguenti:

- *Proprietà degli URL*: alcune analisi sulle proprietà dei link mostrano che gli URL di un host sono delle buone feature per identificare lo spam. In particolare l'URL di un host con molti caratteri, punti, slash e numeri è un buon indicatore

di spam. Un modo semplice per classificare le pagine è quindi quello di usare un valore di soglia che definisca il numero massimo di tali caratteri per identificare una pagina come spam.

- *Host name resolution*: gli spammer possono popolare gli URL delle pagine spam con termini contenuti in query molto frequenti, che sono rilevanti per un certo settore, e impostare un DNS per risolvere questi host name. Gli spammer cercano di manipolare il meccanismo usato da alcuni motori di ricerca (ad esempio Google) che data una query q , assegnano un rank più alto a un URL u se i termini che compongono il nome dell'host di u combaciano con i termini della query. Perciò gli spammer creano tanti URL rilevanti per le diverse query in modo tale da far risultare la pagina di web spam tra i primi risultati di ricerca per molte query. Per determinare tale forma di spam, uno spam detector dovrebbe controllare quanti URL vengono risolti da uno stesso indirizzo IP.
- *Proprietà del contenuto*: le pagine generate automaticamente hanno tutte lo stesso template, ad esempio numerosi siti di spam generano dinamicamente pagine con uno stesso numero di parole. Una tecnica per determinare lo spam è quella di clusterizzare le pagine in base alla somiglianza dei template. Considerando che le pagine di spam hanno strutture molto simili tra loro, se si identificano gruppi con molte pagine aventi la stessa struttura è probabile che esse siano spam.

Oltre a queste proprietà base, in [17] vengono descritti ulteriori metodi e proprietà per l'individuazione dello spam. In tale studio i metodi e le proprietà descritti sono il frutto di analisi effettuate dagli autori su un dataset di pagine HTML, che è stato ricavato utilizzando MSN Search crawler nell'agosto del 2004. Dalle analisi su tale dataset risulta che i domini con maggiore contenuto di spam sono: “.biz”, “.us” e “.com” mentre le pagine contenenti più spam sono: francesi, tedesche e inglesi. I risultati sono rappresentati nei due grafici in figura 2.1 e in figura 2.2. Nel primo grafico l'asse orizzontale rappresenta i domini e l'asse verticale mostra la frazione di spam all'interno di un dominio. I valori nel grafico sono riportati con un intervallo

di confidenza del 95% rappresentato dalla linea verticale sopra ogni barra. L'intervallo di confidenza varia in dimensione a causa del numero differente di campioni raffigurati nei diversi domini. Nel secondo grafico l'asse orizzontale rappresenta la lingua della pagina e l'asse verticale rappresenta la frazione di spam per le pagine scritte in una particolare lingua. Anche nel secondo grafico i valori sono riportati con un intervallo di confidenza del 95%.

Una proprietà che consente di identificare una pagina spam è il numero di parole all'interno della pagina stessa. Infatti una pratica molto comune nel costruire pagine spam è la cosiddetta “Keyword stuffing”, un processo per cui al contenuto della pagina vengono aggiunte molte parole popolari che sono, tuttavia, irrilevanti rispetto ai contenuti; lo scopo è di incrementare le probabilità della pagina di essere in cima ai risultati di molteplici query. In figura 2.3 viene plottato la distribuzione del numero di parole per ogni pagina del dataset (rappresentata dall'istogramma blu) correlata con la probabilità che una pagina sia spam (rappresentata dalla linea viola). Dal grafico si nota che la prevalenza di spam è più alta per le pagine contenenti molte parole. Perciò c'è una correlazione tra prevalenza di spam e numero di parole. Ma si può notare anche che il conteggio delle parole da solo non è una buona euristica visto che porta un alto tasso di falsi positivi.

La tecnica “Keyword stuffing” viene utilizzata anche per la scelta dei titoli, tenendo in considerazione che alcuni motori di ricerca assegnano un peso maggiore ai termini della query presenti all'interno del titolo della pagina. Il grafico in figura 2.4 rappresenta la distribuzione del numero di parole all'interno dei titoli delle pagine correlata con la probabilità che una pagina sia spam. Dal grafico si evince che un eccesso di parole all'interno del titolo è un indicatore di spam. Le parole che vengono utilizzate nel processo di “keyword stuffing” vengono selezionate casualmente o da un ristretto gruppo di query comuni. Per esaminare il comportamento con cui sono selezionate e costruite le frasi innanzitutto vengono identificate le n parole più comuni all'interno del corpus; successivamente viene calcolata, per ogni pagina, la frazione delle parole comuni. Tale processo viene ripetuto per ogni scelta di n (in figura 2.5 $n=200$). Il grafico rappresenta la frazione di parole comuni per ogni pagina; esso ha una caratteristica gaussiana e suggerisce che la maggior parte delle

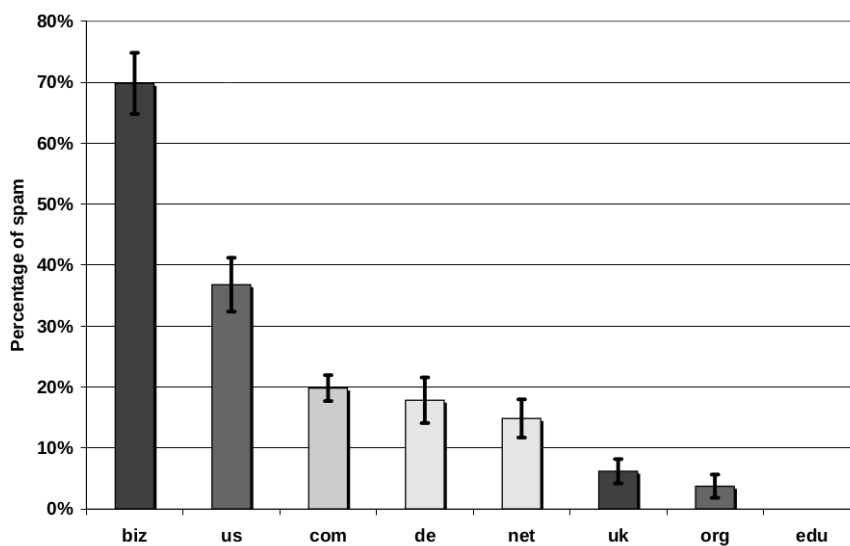


Figura 2.1: Occorrenze dello spam classificate per dominio all'interno del dataset descritto in [17]

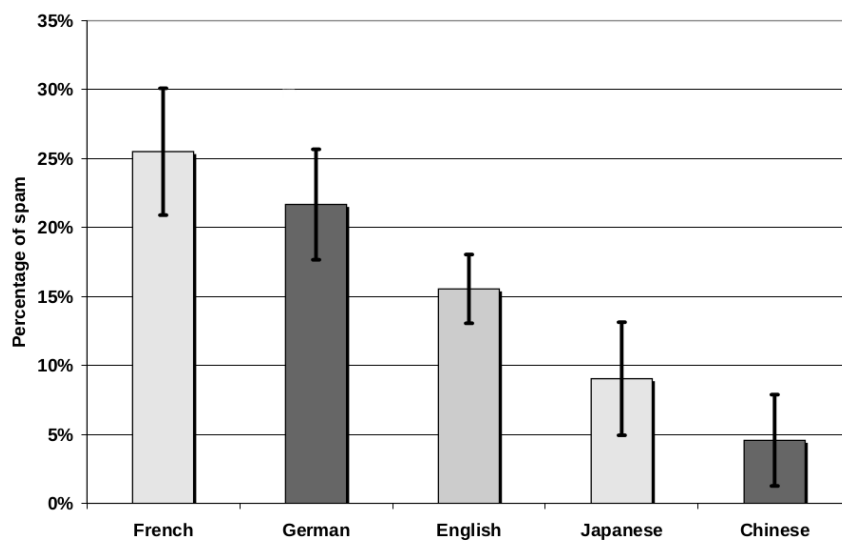


Figura 2.2: Occorrenze dello spam classificate per lingua all'interno del dataset descritto in [17]

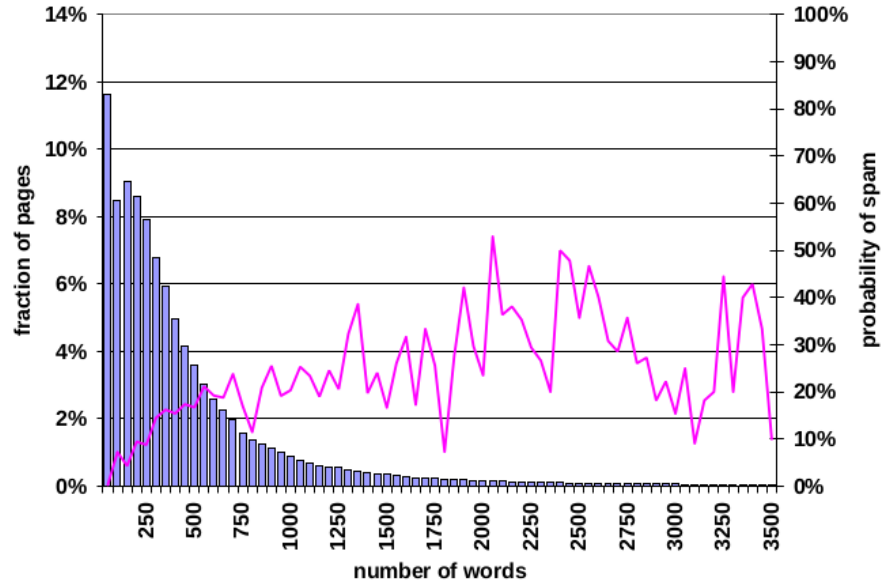


Figura 2.3: Prevalenza di spam sulla base del numero di parole per pagina

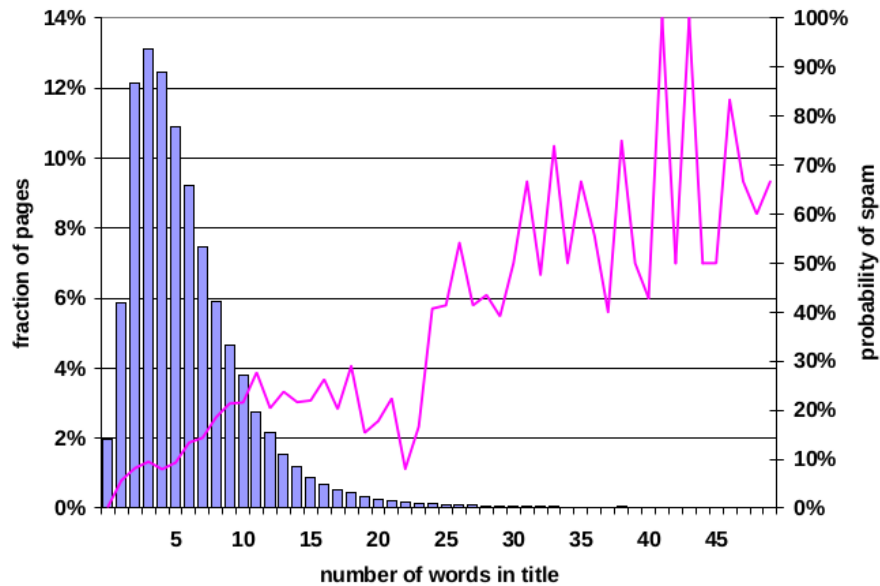


Figura 2.4: Prevalenza di spam sulla base del numero di parole all'interno dei titoli delle pagine

pagine di spam sono generate tessendo parole da un dizionario con una scelta casuale. Un metodo che può essere adottato per identificare le pagine spam che sono generate

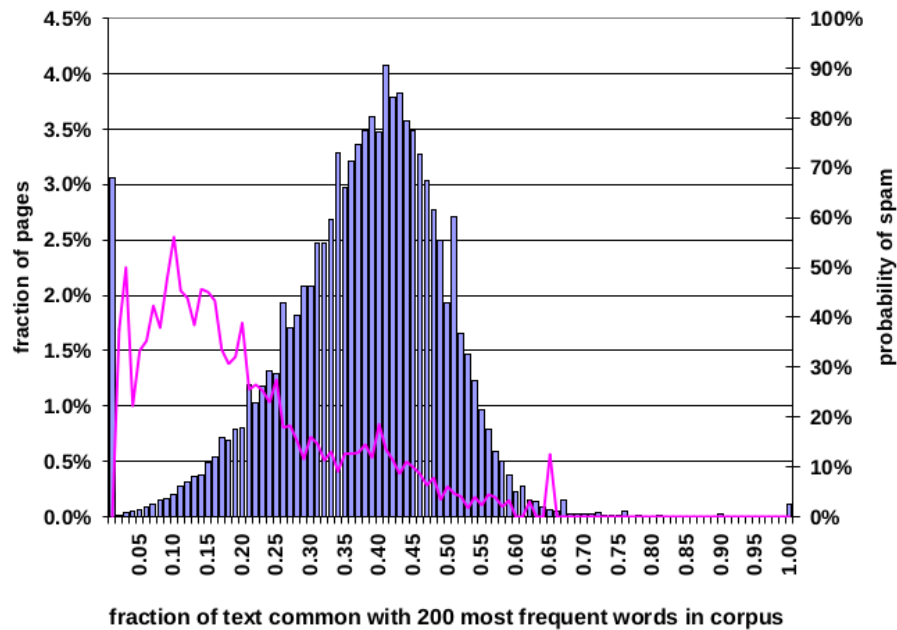


Figura 2.5: Prevalenza di spam sulla base della frazione di parole di una pagina che sono tra le 200 più frequenti parole nel corpus

automaticamente è descritto in [6].

Un altro dato utilizzato per determinare se una pagina è spam è la lunghezza media delle parole delle pagine. Dal dataset preso in considerazione in [17] (grafico in figura 2.6) si nota che la distribuzione della lunghezza media delle parole è simile a una gaussiana con moda e mediana corrispondenti a 5.0 e le parole con lunghezza media uguale a 10 sono certamente spam.

Un'altra proprietà delle pagine web che consente di stimare se una pagina è spam è la quantità di testo che è contenuta all'interno delle ancore (il tag “< a >” delle pagine web). Infatti una pratica comune dei motori di ricerca è considerare il testo delle ancore dei link in una pagina come annotazioni che descrivono il contenuto della pagina che viene puntata. L'idea principale è che se la pagina a ha un link alla pagina b con testo dell'ancora, ad esempio, “computer” allora potremmo concludere che b parli di computer, anche se questa keyword non compare all'interno della

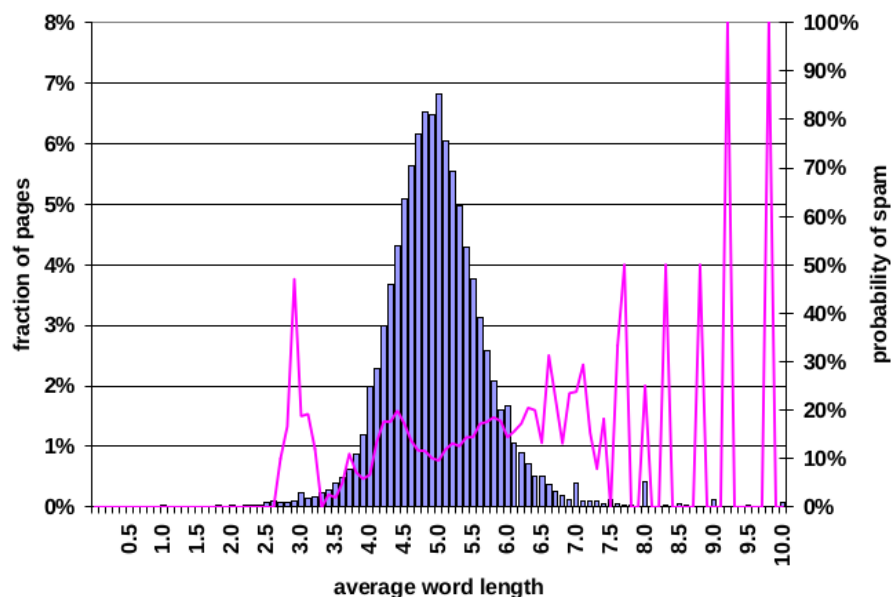


Figura 2.6: Prevalenza di spam sulla base della lunghezza media delle parole per pagina

pagina b . Pertanto durante il ranking tali motori di ricerca potrebbero considerare la pagina b come risultato di una query contenente la keyword “computer”. Sfruttando questo meccanismo si creano pagine di spam contenenti solo del testo all’interno delle ancore per valorizzare il ranking di altre pagine.; queste pagine di norma sono solo cataloghi di link ad altre pagine. Per capire meglio il fenomeno è stato calcolata la frazione di tutte le parole del testo delle ancore all’interno di una pagina, escludendo i markup rispetto al contenuto della pagina. In figura 2.7 viene visualizzato il grafico risultante. Si nota che un’alta frazione di testo delle ancore aumenta la probabilità che la pagina sia spam ma usare questa euristica da sola potrebbe portare un alto numero di falsi positivi [17].

Calcolando la frazione di contenuto visibile all’interno di una pagina (definita come la lunghezza in termini di byte di tutte le parole non di markup) rispetto all’intera dimensione della pagina si nota dal grafico, in figura 2.8, che la distribuzione (delle frazioni di contenuto visibile) evidenzia che le pagine di spam hanno meno markup delle pagine normali. Questo fa intendere che molte pagine spam hanno il solo scopo di dover essere indicizzate dai motori di ricerca e non di essere fruite da un utente.

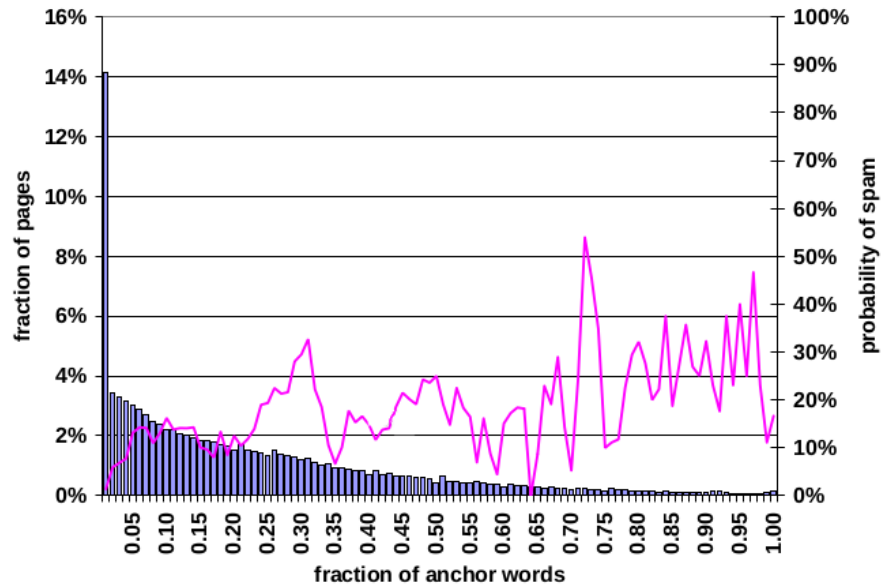


Figura 2.7: Prevalenza di spam sulla base della quantità di testo delle ancore delle pagine

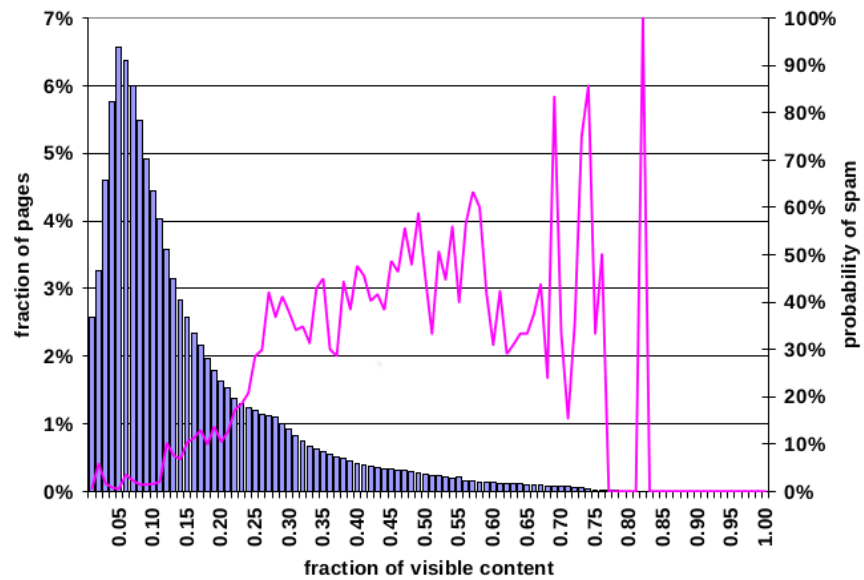


Figura 2.8: Prevalenza di spam sulla base della frazione di contenuto visibile

Come detto in precedenza i motori di ricerca possono dare un peso maggiore a pagine che contengono ripetutamente le keyword contenute nella query (ad esempio utilizzano come metodo di ranking il “term-frequency”). Le pagine spam hanno, quindi, contenuti replicati molte volte per aumentare il rank. Per rilevare la ridondanza di contenuti (ottenuta dal processo di replica), viene calcolato il rapporto di compressione ovvero la dimensione della pagina non compressa divisa per la dimensione della pagina compressa. In figura 2.9 è rappresentata la distribuzione del rapporto di compressione e la likelihood che la pagina sia spam. Dal grafico si nota che quanto più il valore di compressione è elevato tanto più probabilmente la pagina può essere considerata spam; il 70 per cento delle pagine con un rapporto di compressione maggiore di 4.0 sono giudicate spam. Questo è dovuto al fatto che una compressione su una pagina spam che ha contenuti ridondanti, sarà più efficace di una compressione su una pagina non spam che è caratterizzata da contenuti non ridondanti. E perciò il rapporto di compressione tenderà a crescere quanto più la pagina sarà caratterizzata da contenuti ridondanti.

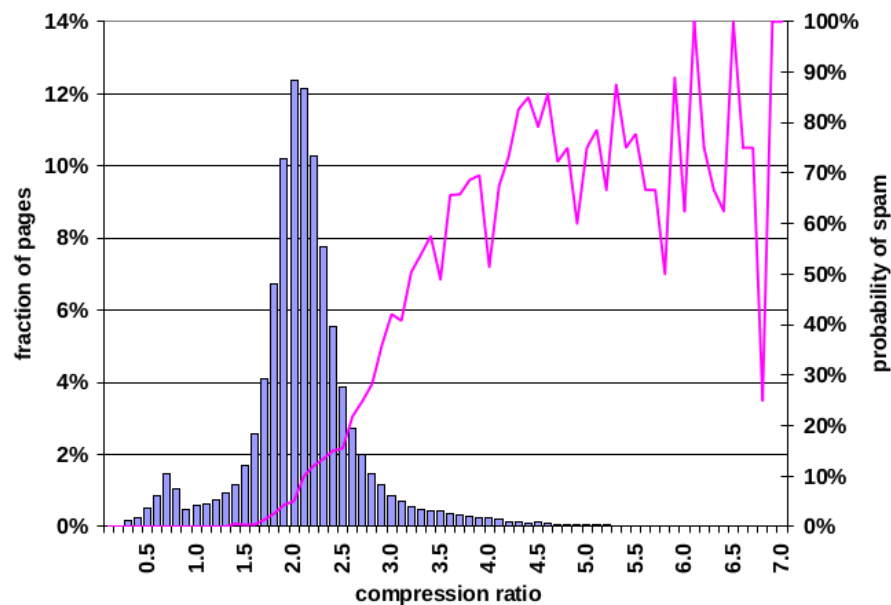


Figura 2.9: Prevalenza di spam sulla base del rapporto di compressione

2.1.2 Utilizzo di un classificatore per combinare le feature

Sempre in [17] le euristiche descritte fino a questo punto sono combinate considerando il problema di rilevamento dello spam come un problema di classificazione. Quindi viene creato un classificatore, che usando le caratteristiche di una pagina web la classificherà in una delle due classi: spam o non spam. Costruire un classificatore richiede una fase di training durante la quale i parametri del classificatore sono determinati e una fase di testing durante la quale le performance del classificatore sono valutate. Per ogni pagina all'interno del dataset viene calcolato il valore di ogni feature (ad esempio, usando le euristiche discusse sopra) e si utilizzano questi valori per istruire il classificatore. Il classificatore usato è un albero di decisione che dato un insieme di dati da training e un insieme di feature crea un diagramma di flusso ad albero. Ogni nodo dell'albero corrisponde al test da valutare per una particolare feature mentre ogni arco è un valore di uscita del test ed infine le foglie corrispondono alle classi che verranno assegnate alle pagine. Un esempio del classificatore è rappresentato in figura 2.10. Per migliorare il classificatore si possono usare tecniche come bagging o boosting. Queste tecniche creano un insieme di classificatori che vengono combinati.

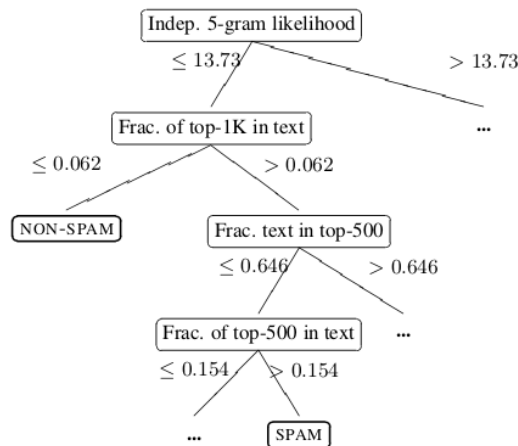


Figura 2.10: Esempio di classificatore

2.1.3 Language model per rilevare lo spam

Oltre alle feature presentate in precedenza, in [16] sono descritte nuove tipologie di feature e metodi per rilevare lo spam che fanno uso dei *language model* per analizzare le sorgenti estratte da ogni sito in un dataset. Le sorgenti sono dei pezzi di contenuto presenti all'interno delle pagine web di partenza (quelle che contengono i link ad altre pagine) e all'interno delle pagine web di destinazione (quelle che sono linkate dalle pagine di partenza). Per ogni sorgente si definisce un modello e si calcola quanto siano differenti i modelli individuati; viene utilizzata la *Kullback-Leibler Divergence* (*KLD*) per misurare la divergenze tra le distribuzioni di probabilità dei termini di pagine web, applicandola a unità di testo della pagina di partenza e di quella linkata. La Kullback-Leibler (KL) è una misura asimmetrica della divergenza che misura quanto male una distribuzione di probabilità M_q riesce a modellare M_d :

$$KLD(T_1||T_2) = \sum_{t \in T_1} P_{T_1}(t) \log \frac{P_{T_1}(t)}{P_{T_2}(t)} \quad (2.1)$$

dove in 2.1 $P_{T_1}(t)$ è la probabilità del termine t nella prima unità di testo e $P_{T_2}(t)$ è la probabilità del termine t nella seconda unità di testo. In basso sono rappresentati due esempi di KLD applicata tra il testo delle ancore della pagina sorgente e i titoli delle pagine puntate dai link (esempio preso da WEBSHAM-UK2006). Dagli esempi si deduce che nel primo caso c'è una maggiore correlazione tra il testo dell'ancora della pagina sorgente e il titolo della pagina puntata rispetto al secondo esempio dove il valore di KLD è maggiore.

```
KLD(Free Ringtones || Free Ringtones for Your Mobile Phone from
PremieRingtones.com) = 0.25

KLD(Best UK Reviews || Findabmw.co.uk - BMW Information
Resource) = 3.89
```

Per determinare se una pagina è spam si cerca di identificare una relazione tra due pagine collegate sulla base del valore di divergenza. I valori sono ottenuti calcolando le divergenze con KLD tra una o più sorgenti di informazioni da ogni pagina. In particolare si usano tre tipi di informazione di una pagina sorgente: testo delle

ancore, testo intorno alle ancore, termini nell'URL mentre per la pagina che viene linkata dalla pagina sorgente: titolo, contenuto della pagina e meta tag. Combinando queste sorgenti di informazione si determina la divergenza tra due pagine; di seguito vengono descritte alcune combinazioni base:

- testo delle ancore - contenuto. Una pagina che ha un link verso un'altra pagina ha solo un modo per convincere l'utente a visitare la pagina collegata: mostrare in maniera concisa le informazioni relative alla pagina collegata. Perciò una grande divergenza tra questi due frammenti di testo indica che la pagina può essere spam. In figura 2.11 è illustrata la divergenza KL tra le due sorgenti di informazione, come si vede la curva delle pagine normali è più compatta di quella delle pagine spam.

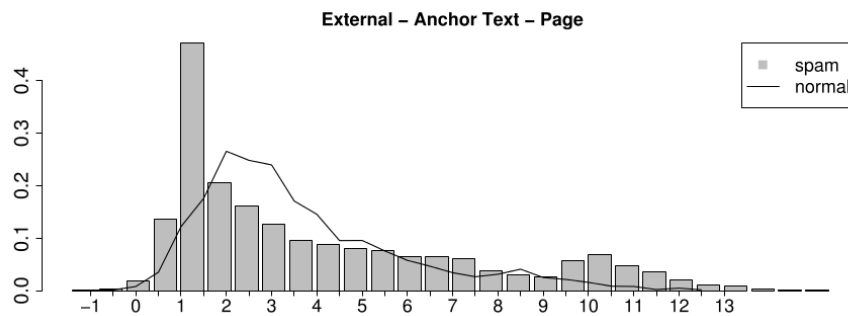


Figura 2.11: Istogramma della divergenza KL tra testo delle ancore e il contenuto della pagina puntata basato sul dataset *WEBSPAM-UK2006* utilizzato in [16]

- testo vicino alle ancore - contenuto. Il testo delle ancore, a volte, è un valore poco descrittivo e per ovviare al problema viene usato il testo che circonda le ancore. Nell'esperimento vengono utilizzate 7 parole per lato. In figura 2.12 viene mostrato che le pagine spam hanno alti valori di divergenza mentre le normali sono concentrate intorno $KL \approx 2.5$, .
- termini nell'URL - contenuto. I motori di ricerca danno molta importanza agli URL perché i termini che li compongono possono dare un'idea del contenuto della pagina a cui l'URL fa riferimento. Un metodo di spam che sfrutta questo meccanismo consiste nella creazione di URL contenenti molti termini comu-

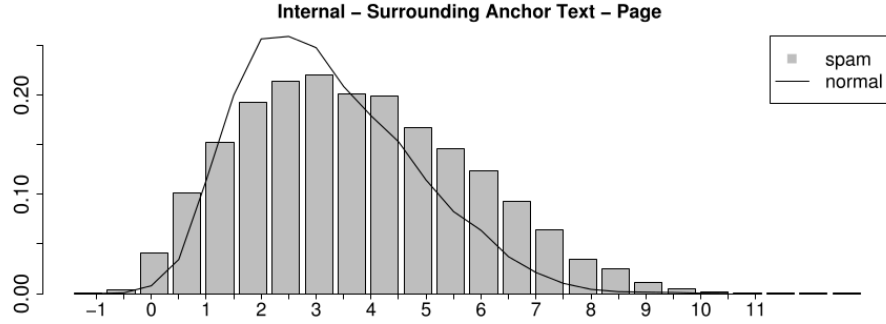


Figura 2.12: Istogramma della divergenza KL tra testo intorno alle ancore e il contenuto della pagina puntata basato sul dataset WEBSPPAM-UK2006 utilizzato in [16]

ni (ad esempio: `www.domain.com/viagra-youtube-free-download-poker-online.html` che in realtà è solo uno store online e non ha alcuna attinenza con i nessuno dei termini utilizzati nell'URL). Al fine di identificare questo tipo di spam, vengono prelevati i termini più rilevanti da un URL e calcolata la divergenza col contenuto della pagina di destinazione. La distribuzione finale, rappresentata in figura 2.13, dimostra una grande differenza tra l'istogramma delle pagine normali con quello delle pagine spam.

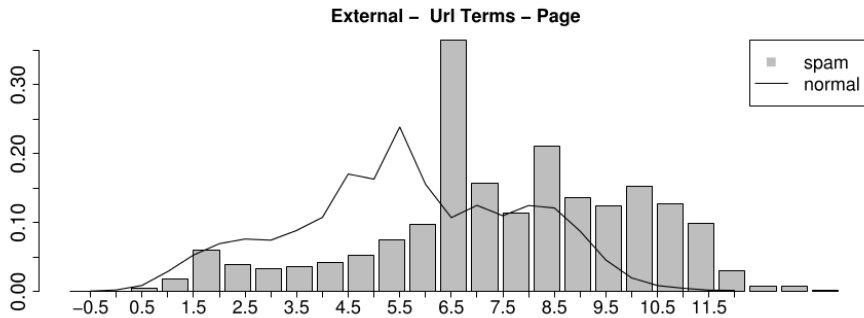


Figura 2.13: Istogramma della divergenza KL tra termini degli URL e il contenuto della pagina puntata basato sul dataset WEBSPPAM-UK2007 utilizzato in [16]

- testo delle ancore - titolo. Questa feature si compone di due sorgenti che sono molto simili in quanto descrivono la pagina con poche parole. Tuttavia la prima sorgente può essere scritta anche da chi non è il proprietario della pagina di destinazione. In figura 2.14 notiamo che questa feature da sola non discrimina

bene le pagine spam ma è abbastanza efficace se utilizzata in congiunzione con altre feature.

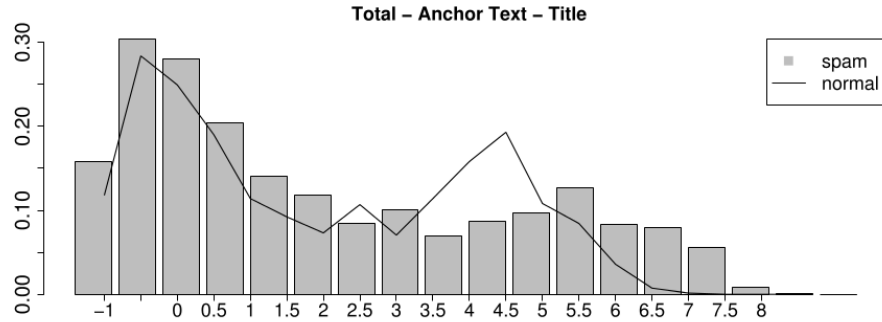


Figura 2.14: Istogramma della divergenza KL tra testo delle ancore e titolo della pagina puntata basato sul dataset *WEBSPAM-UK2007* utilizzato in [16]

- testo intorno alle ancore - titolo. Dal grafico in figura 2.15 si può notare come tale feature rileva meglio lo spam rispetto alla precedente, infatti molti valori di spam sono concentrati per valori di $KL > 3$.

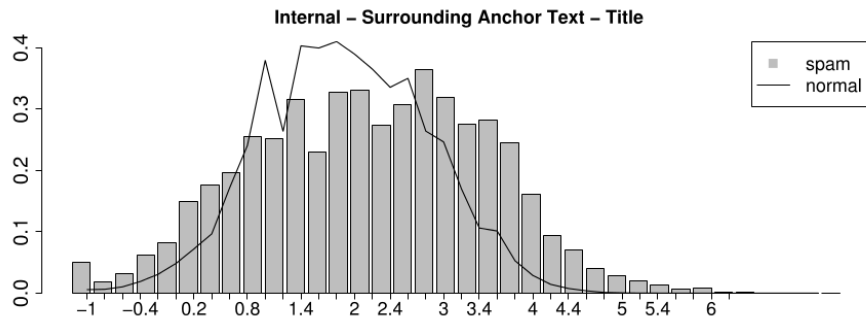


Figura 2.15: Istogramma della divergenza KL tra testo intorno alle ancore e titolo della pagina puntata basato sul dataset *WEBSPAM-UK2006* utilizzato in [16]

- termini nell'URL - titolo. In questa feature entrambe le sorgenti sono generate dal proprietario della pagina, a differenza della precedente in cui la sorgente di informazione della pagina di partenza poteva essere generata da un'altra persona. In questo vi dovrebbe essere maggiore coerenza tra le due sorgenti. In figura 2.16 sono rappresentate le distribuzioni per le pagine spam e non spam.

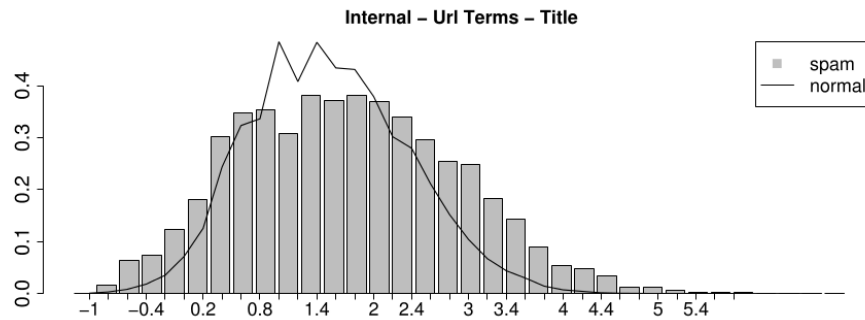


Figura 2.16: Istogramma della divergenza KL tra termini nell'URL e titolo della pagina puntata basato sul dataset WEBSPAM-UK2006 utilizzato in [16]

- titolo - contenuto. I motori di ricerca danno un peso maggiore ai termini della query se presenti nel titolo della pagina. Sfruttando tale meccanismo gli spammer perfezionano i loro processi in modo tale da impostare termini chiave nel titolo anche creando però una divergenza tra titolo e contenuto della pagina. La feature titolo - contenuto consente la rilevazione di spam quando non vi è alcuna relazione tra il titolo e il contenuto della pagina. In figura 2.17 è rappresentata la divergenza tra le due distribuzioni spam e non spam.

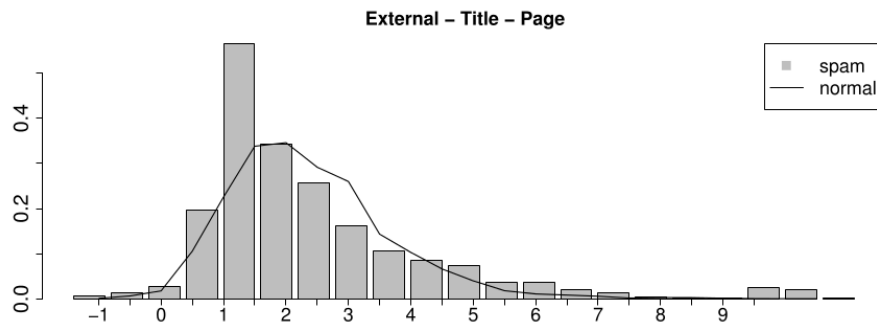


Figura 2.17: Istogramma della divergenza KL tra titolo e contenuto della pagina basato sul dataset WEBSPAM-UK2006 utilizzato in [16]

- metatag. Vengono usati per calcolarne la divergenza con altre sorgenti di informazioni della pagina di partenza (come il testo delle ancore e il testo intorno alle ancore) e della pagina destinazione (come il contenuto o i termini

dell'URL). In figura 2.18 viene visualizzata la divergenza tra il testo delle ancore e i metatag.

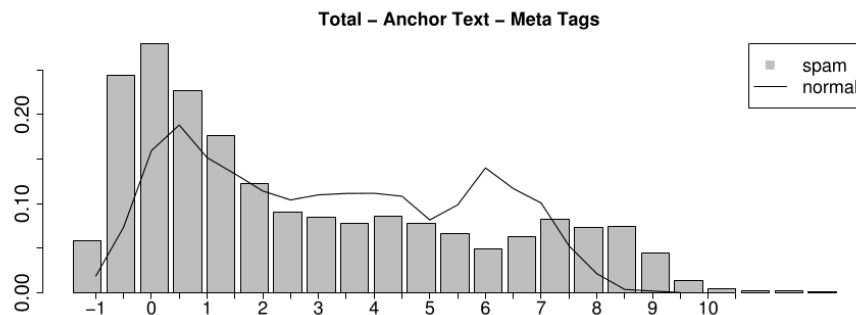


Figura 2.18: Istogramma della divergenza KL tra testo delle ancore e i meta tag della pagina basato sul dataset *WEBSPAM-UK2006* utilizzato in [16]

Oltre alle feature descritte si possono ottenere delle feature più articolate combinando le feature base tra di loro; ad esempio per le pagine sorgenti si possono definire: testo delle ancore e termini nell'URL, testo intorno alle ancore e URL (per maggiori dettagli [16]).

Infine tali feature possono essere usate per istruire un classificatore.

2.1.4 Spam detection sulla base degli argomenti di una pagina web

Dong et al. [3] propongono un metodo basato su statistiche effettuate sugli argomenti delle pagine. Le statistiche non si basano sulle parole della pagina ma sugli argomenti, in modo da non ignorare la semantica delle parole e di catturare le feature linguistiche nascoste nel testo per capire se la pagina è spam. Le analisi vengono fatte usando i topic model, modelli statistici che scoprono gli argomenti latenti presenti in una collezione di documenti. Un topic model utilizzato è la *Latent Dirichlet Allocation (LDA)* che modella ogni argomento latente come una distribuzione probabilistica su un vocabolario e ogni documento come una distribuzione probabilistica sugli argomenti latenti. L'intuizione di usare i topic model per analizzare il contenuto delle pagine web nasce dal fatto che analizzando le feature nascoste di una pagina spam, gli autori hanno notato che i contenuti di queste pagine, generati automati-

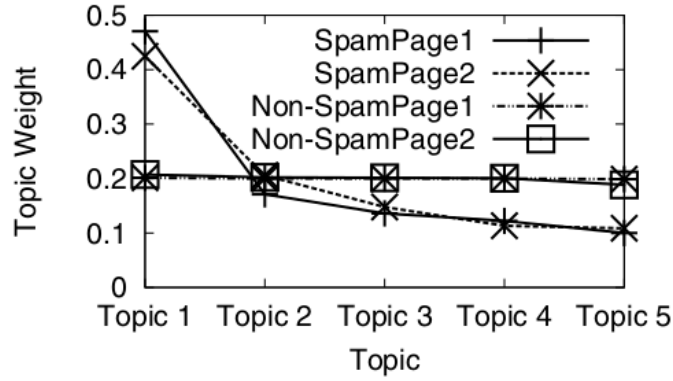


Figura 2.19: Distribuzione degli argomenti pesati per pagine spam e normali

camente, sono differenti dai contenuti delle pagine non spam. Vengono definiti tre tipi di misure.

La prima misura per determinare le pagine spam utilizzando LDA sfrutta la caratteristica che tali pagine sono molto topic-centric, ovvero hanno uno specifico insieme di argomenti. In figura 2.19 sono rappresentate quattro distribuzioni degli argomenti presenti in quattro pagine (due spam e due non spam) scelte casualmente da un dataset. A supporto della tesi degli autori che le pagine spam sono topic centric si può notare che presentano una distribuzione esponenziale, in quanto esse vengono sviluppate per avere un alto ranking per un insieme di specifiche query di ricerca, al contrario delle pagine non spam che presentano una distribuzione uniforme (nell'articolo gli autori sostengono che le pagine non spam come una homepage contengono vari argomenti come ad esempio: i contatti, chi e cosa fa, altre informazioni). Per classificare le pagine sulla base di questa caratteristica, basata sulle distribuzioni degli argomenti, è stato proposto una misura della diversità degli argomenti basata sulla varianza. Data una pagina web d , la sua distribuzione degli argomenti è $T(d) = \{t_1, t_2, \dots, t_m\}$, dove ogni argomento $t_i (1 \leq i \leq m)$ è associato con un peso δ_{t_i} . La misura della diversità degli argomenti basata sulla varianza per d , denotata con $TopicVar(d)$, è calcolata come:

$$TopicVar(d) = \frac{\sum_{i=1}^m (\delta_{t_i} - u)^2}{m} \quad (2.2)$$

dove $u = \frac{\sum_{i=1}^m \delta_{t_i}}{m} = \frac{1}{m}$. In figura 2.20 è illustrata la distribuzione risultante. Dalla

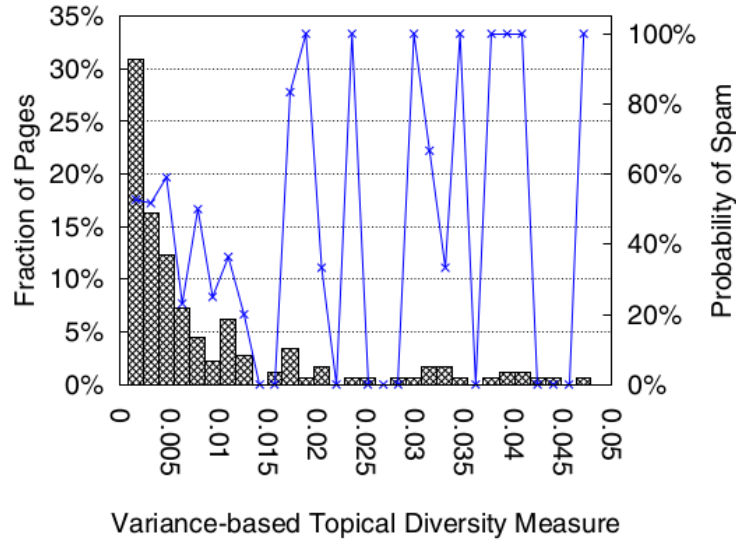


Figura 2.20: Prevalenza di spam relativa alla misura della diversità degli argomenti basata sulla varianza

distribuzione si nota che le pagine spam sono più topic-centric ovvero la varianza dei pesi degli argomenti è più grande rispetto alle pagine non spam. Questo è dovuto al fatto che le pagine non spam avendo una distribuzione quasi uniforme hanno una la varianza che è più piccola rispetto alle pagine spam che hanno una distribuzione più concentrata di quella uniforme, come quella esponenziale. Il valore di *TopicVar* è proporzionale alla probabilità che una pagina sia spam, ovvero all'aumento della probabilità di una pagina di essere spam consegue un aumento del valore di *TopicVar* per quella pagina. Questa misura è un ottimo indicatore per rilevare lo spam.

La seconda misura per identificare le pagine spam utilizzando LDA che permette di misurare la relazione semantica tra gli argomenti è la semantica delle parole. Gli autori definiscono che: data una pagina web d , la sua distribuzione degli argomenti è $T(d) = \{t_1, t_2, \dots, t_m\}$. La probabilità che una parola w appartenga a un argomento t_i ($1 \leq i \leq m$) è definita come $\phi(w|t_i)$. Ogni argomento t_i è rappresentato come un insieme di parole denotate come $W(t_i)$. Intuitivamente due argomenti t_i, t_j ($1 \leq i, j \leq m$) sono semanticamente correlati se le parole $W(t_i)$ e $W(t_j)$ sono semanticamente legate. In [3] per ottenere le relazioni semantiche tra le due parole è utilizzata una funzione di similarità $Sim(w_i, w_j)$ che è quella di Wordnet. Quindi

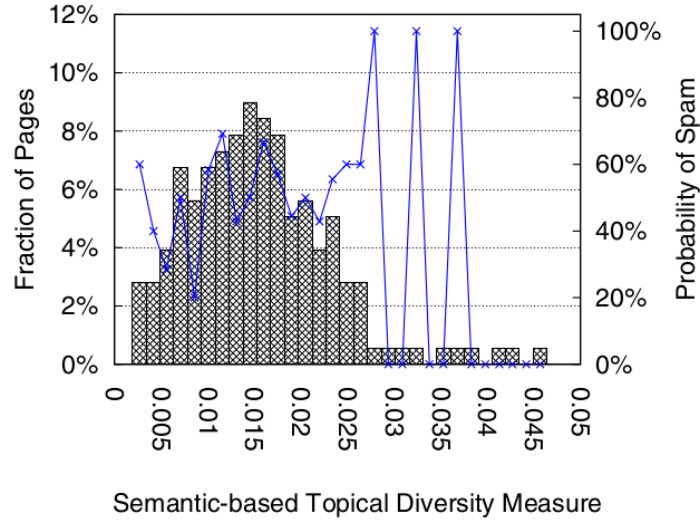


Figura 2.21: Prevalenza di spam relativa alla misura di diversità degli argomenti basata sulla semantica

per misurare la relazione semantica tra due argomenti t_i, t_j , si calcolano le similarità tra ogni coppia di parole dei due argomenti moltiplicate con le loro probabilità rispetto agli argomenti.

$$Sim(t_i, t_j) = \frac{\sum_{w_k \in W(t_i), w_l \in W(t_j)} Sim(w_k, w_l) X\phi(w_k|t_i) X\phi(w_l|t_j)}{\frac{|W(t_i)|X|W(t_j)|}{2}} \quad (2.3)$$

Usando quindi un modello degli argomenti otteniamo m argomenti latenti. Dall'equazione 2.3 si deriva una misura della diversità degli argomenti basata sulla semantica per gli m argomenti latenti di una collezione [3]: data una pagina web d , la sua distribuzione degli argomenti è $T(d)$ e quindi la misura della diversità degli argomenti basata sulla semantica per tale pagina d è:

$$TopicSim(d) = \frac{\sum_{1 \leq i \leq j \leq m} Sim(t_i, t_j)}{\frac{1}{2}m(m-1)} \quad (2.4)$$

In figura 2.21 è illustrata la distribuzione della misura di diversità degli argomenti basata sulla semantica. Dalla distribuzione si nota che quando la misura cresce la probabilità che una pagina sia spam aumenta, ovvero le pagine spam hanno argomenti che sono in forte relazione semantica tra loro e cioè sono icentrate su pochi argomenti.

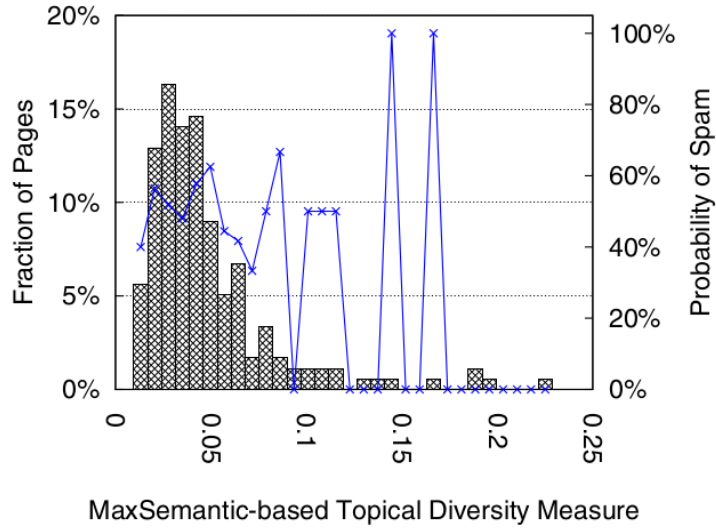


Figura 2.22: Prevalenza di spam relativa alla misura di diversità sulla massima semantica

L'ultima misura è basata sulla massima semantica ed è definita come:

$$TopicSimMax(d) = \max\{Sim(t_i, t_j) | 1 \leq i \leq j \leq m\} \quad (2.5)$$

La distribuzione della misura di diversità basata sulla massima semantica mostra che questa è più alta per le pagine che contengono spam. La distribuzione è rappresentata in figura 2.22.

Anche in questo caso per determinare se una pagina è spam oppure non spam, vengono utilizzati algoritmi supervisionati di apprendimento per istruire un classificatore di pagine spam usando misure di diversità degli argomenti. Con LDA possiamo impostare parametri come il numero di argomenti e il numero di parole per ogni argomento che incidono sulle le prestazioni della classificazione.

2.1.5 Altre tecniche

Nello studio in [1] viene presentato l'algoritmo WITCH (Web Spam Identification Through Content and Hyperlinks), un algoritmo ibrido che utilizza sia il contenuto della pagina che la struttura dei link per identificare le pagine spam. Come descritto in precedenza nel sotto capitolo 2.1.1, le differenti proprietà tra le pagine spam e non spam possono essere sfruttate per costruire un classificatore. Per identificare lo spam l'algoritmo WITCH utilizza le feature descritte in precedenza e analizza la struttura

dei collegamenti tra le pagine. In particolare viene istruito un classificatore lineare nello spazio delle feature usando la SVM (Support Vector Machine) come funzione obbiettivo. I collegamenti tra le pagine sono utilizzati in modo da regolarizzare il grafo, che produce una predizione che varia leggermente tra le pagine dei link. Il metodo SVM associato alla regolarizzazione del grafo è efficiente per il rilevamento di web spam.

Un altro metodo, proposto in [21], è denominato “Hidden style similarity measure” ed è basato su feature extra testuali appartenenti alle pagine HTML. Infatti gli autori sostengono che le pagine spam generate automaticamente non sono facili da rilevare utilizzando metodi classici di classificazione basati solamente sul contenuto; gli autori quindi utilizzano la struttura HTML di una pagina per classificare le pagine simili.

2.2 Tecniche basate sul grafo

Tali tecniche fanno uso del grafo del web ricavato dai collegamenti ipertestuali tra le pagine. Il web, quindi, può essere rappresentato come un grafo diretto $G = (V, E)$, dove V è l'insieme delle pagine e rappresentano i nodi del grafo mentre E è l'insieme dei link diretti tra le pagine. Il grafo può essere astratto e rappresentato da una matrice di transizione così formata:

$$T(p, q) = \begin{cases} 0 & \text{if } (q, p) \in E \\ 1/\omega(q) & \text{if } (q, p) \in E \end{cases} \quad (2.6)$$

dove $\omega(p)$ è il grado di link in uscita della pagina p . Possiamo anche definire la matrice di transizione inversa U :

$$U(p, q) = \begin{cases} 0 & \text{if } (p, q) \in E \\ 1/l(q) & \text{if } (p, q) \in E \end{cases} \quad (2.7)$$

dove $l(q)$ è il grado di link in ingresso della pagina q .

2.2.1 Metodi classici per identificare lo spam web usando il grafo

Uno dei primi metodi adottati per identificare lo spam web usando il grafo è *Trustrank* [9]. *Trustrank* fa uso di un insieme di pagine di partenza S che sono valutate

da degli esperti e che vengono classificate in due sottoinsiemi: pagine non spam S^+ e pagine spam S^- ; questa fase è chiamata funzione *Oracle*. Per determinare le pagine non spam senza invocare la funzione *Oracle* su tutto il grafo derivato dalla fase di crawling, viene fatta un'assunzione empirica chiamata *isolazione approssimata dell'insieme delle pagine buone* la quale afferma che le pagine non spam raramente punteranno a quelle spam perché gli sviluppatori di pagine non spam hanno poco interesse nel linkare pagine spam almeno che non vengano ingannati tramite ad esempio l'uso di tecniche come l'*honeypot*. Dato un numero limitato di chiamate della funzione *Oracle* sul seed set di partenza e sfruttando l'assunzione fatta precedentemente viene definita una funzione, denominata come *funzione di verità ignorante* T_0 , per ogni pagina p del grafo:

$$T_0(p) = \begin{cases} O(p) & \text{if } p \in S \\ 1/2 & \text{altrimenti} \end{cases} \quad (2.8)$$

dove la funzione O è la funzione *Oracle*. Dal momento che le pagine buone dovrebbero puntare ad altre pagine buone assegnamo 1 a tutte le pagine che possono essere raggiunte da una pagina in S^+ in M step. La funzione di verità T_M è definita come:

$$T_M(p) = \begin{cases} O(p) & \text{if } p \in S \\ 1 & \text{if } p \notin S \text{ and } \exists q \in S^+ : q \rightarrow_M p \\ 1/2 & \text{altrimenti} \end{cases} \quad (2.9)$$

Il percorso dalla pagina q a p nell'equazione non comprende pagine spam incluse nell'insieme S^- .

Il problema della funzione di verità T_M è che non esiste la sicurezza che le pagine raggiungibili da pagine buone siano effettivamente della stessa caratteristica. Infatti più lontana una pagina p si trova dal seed set S^+ minore è la certezza che quella pagina sia buona. Un modo per non incorrere in questo errore è ridurre il valore della funzione di verità ogni qual volta ci si allontana dal seed set S^+ .

In figura 2.23 è possibile vedere in dettaglio l'algoritmo. L'algoritmo calcola il valore di verità di ogni pagina dell'intero grafo. I valori di input sono il grafo descritto dalla matrice di transizione T e il numero di pagine N e i parametri di controllo dell'esecuzione: L il numero di chiamate della funzione *Oracle* e α_b il fattore di

```

function TrustRank
input
    T      transition matrix
    N      number of pages
    L      limit of oracle invocations
     $\alpha_B$  decay factor for biased PageRank
     $M_B$     number of biased PageRank iterations
output
    t*     TrustRank scores
begin
    // evaluate seed-desirability of pages
(1)  s = SelectSeed(...)
    // generate corresponding ordering
(2)   $\sigma = \text{Rank}(\{1, \dots, N\}, \mathbf{s})$ 
    // select good seeds
(3)  d =  $\mathbf{0}_N$ 
    for  $i = 1$  to L do
        if  $O(\sigma(i)) == 1$  then
            d( $\sigma(i)$ ) = 1
    // normalize static score distribution vector
(4)  d = d / |d|
    // compute TrustRank scores
(5)  t* = d
    for  $i = 1$  to  $M_B$  do
        t* =  $\alpha_B \cdot \mathbf{T} \cdot \mathbf{t}^* + (1 - \alpha_B) \cdot \mathbf{d}$ 
    return t*
end

```

Figura 2.23: Algoritmo di *trustrank*

decadimento per il calcolo di *Pagerank* ed infine M_b il numero di iterazioni per il calcolo di *Pagerank*. Al primo passo viene invocata la funzione *SelectSeed()* calcola l'insieme delle pagine con il relativo rank di rilevanza per essere incluse nel seedset di partenza. Nel secondo punto la funzione *Rank*(x, s) ordina gli elementi di x in modo decrescente sulla base dello score di s . Il punto tre invoca la funzione *Oracle* su L pagine. I valori del vettore d che corrispondono alle pagine buone del seed sono impostate a 1. Nel punto (4) il vettore viene normalizzato in modo tale che la somma faccia 1. Infine al punto (5) viene calcolato *Trustrank* usando *Pagerank* personalizzato dal vettore d che rimpiazza la distribuzione uniforme. Dall'algoritmo si nota che *Trustrank* è una versione modificata di *Pagerank* dove il vettore di teletrasporto è il seed set S^+ calcolato al punto 3 e 4.

Un altro algoritmo che è stato progettato per identificare lo spam usando come

input il grafo delle pagine web è *Anti-Trust Rank* [13]. Questo algoritmo sfrutta la stessa intuizione di *Trustrank* dell'isolamento approssimato cioè che pagine buone molto raramente punteranno a pagine malevoli; quindi si popola un seed set formato da pagine spam e si propaga la funzione Anti Trust (che sarebbe la funzione di verità di *Trustrank*) sul grafo trasposto con l'obbiettivo di rilevare le pagine spam, le quali successivamente possono essere filtrate da un motore di ricerca. Più precisamente a differenza per quanto avviene in *Trustrank* dove la funzione *Trust* è propagata dal seed set composto da pagine non spam lungo tutto il grafo, in *Anti-Trust Rank* la funzione (in questo caso la funzione *Anti Trust*) è propagata nella direzione inversa ai link in entrata ad ogni pagina del grafo, partendo da un insieme di pagine del seed set composto da pagine spam. L'obbiettivo è assegnare un rank maggiore alle pagine spam e successivamente eliminarle dalle ricerche o usando un valore di soglia oppure ritornando le n pagine che hanno valore di *Anti-Trust Rank* più alto.

Trustrank e *Anti-Trust rank* sono ottimi algoritmi per identificare lo spam, ma hanno il problema che l'insieme seed usato potrebbe non essere sufficientemente rappresentativo per coprire bene tutti gli argomenti del web. Un modo naturale di ottenere una grande copertura del web è usare gli argomenti delle pagine come segnale di ingresso: invece di usare un singolo valore di *trustrank* per un sito, in [22] gli autori propongono di calcolare *trustrank* per i differenti argomenti di ogni sito. L'algoritmo consiste nel partizionare il seed set sulla base dei vari argomenti che esso contiene e usare ognuna di queste partizioni come seed set per calcolato il valore di *trustrank* per ogni pagina.

Un altro metodo per l'identificazione di pagine spam è descritto in [2]. Questo metodo separa la credibilità di una pagina dalla credibilità del link per quella pagina al contrario di *pagerank* che è manipolabile tramite tecniche come *hoenypot*. La credibilità viene definita in termini di credibilità *k-scope*. Data una funzione C essere una funzione di credibilità che istantaneamente valuta la qualità di un link di una pagina p al tempo t , un valore di $C(p, t) = 0$ indica che p non è credibile mentre $C(p, t) = 1$ indica che p è credibile. Dato un percorso in un grafo diretto G dalla pagina p alla pagina q essere la sequenza di nodi: $path(p, q) = (n_0, n_1, \dots, n_j)$ dove $p = n_0, q = n_j$ tale che esiste un arco diretto tra nodi successivi nel percorso

$n_i, n_{i+1} \in L$ per $0 \leq i \leq j-1$, diciamo che un percorso in un grafo diretto G dalla pagina p alla pagina q è un *bad path* se la pagina di destinazione è una pagina spam $q \in P_b$ (dove P_b è l'insieme delle pagine spam) e nessuna altra pagina nel percorso è una pagina spam. $path(p, q) = (n_0, n_1, \dots, n_j)$ e $q \in P_b$ e $n_i \notin P_b (0 \leq i \leq j-1)$. La probabilità che una camminata casuale passi, lungo un percorso di lunghezza k , da una pagina p è denotata con $Pr(path_k(p))$ ed è determinata con i pesi degli archi per ogni hop nel percorso:

$$PR(path_k(p)) = \prod_{i=0}^{k-1} w(n, n_{i+1}) \quad (2.10)$$

Quindi credibilità *k-scope* di una pagina è definita in termini di probabilità che una camminata casuale eviti le pagine spam dopo aver superato k hop dalla pagina di origine. La credibilità *k-scope* di una pagina p al tempo t , denotata con $C_k(p, t)$ è definita come segue:

$$C_k(p, t) = 1 - \sum_{l=1}^k \left(\sum_{path_l(p) \in BPath_l(p)} Pr(path_l(p)) \right) \quad (2.11)$$

Nel caso $p \in P_b$ allora $C_k(p, t) = 0$. Nel caso in cui non ci siano pagine spam all'interno di k hop di pagine allora p è credibile con un valore $C_k(p, t) = 1$ se lei è un pagina spam o nel caso in cui tutti i percorsi originati da p colpiscono una pagina p all'interno di k hop, allora p non è credibile $C_k(p, t) = 0$. Ma dato che non è possibile avere tutto il grafo e non c'è nessuna sicurezza sulla conoscenza totale dei nodi spam è stato introdotto il concetto di credibilità tunable k-Scope, la quale aumenta il calcolo della credibilità k-scope includendo un fattore di penalità di credibilità. Gli obiettivi sono approssimare al meglio la credibilità k-scope sotto limiti reali e capire come parametri differenti protrebbero influire sulla qualità delle varie funzioni usate. Sia $G = (P, L)$ essere un grafo diretto, k il raggio massimo di camminata e $\gamma(p)$ il fattore di penalità di credibilità di una pagina $p \in P$ dove $0 \leq \gamma(p) \leq 1$. Definiamo la credibilità tunable k-scope di una pagina p , denotata con $C_k(p)$, in due fasi, quando $p \notin P_b$:

$$C_k(p) = \left(1 - \sum_{l=1}^k \left(\sum_{path_l(p) \in BPath_l(p)} Pr(path_l(p)) \right) \right) \cdot \gamma(p) \quad (2.12)$$

e quando $p \in P_b$ allora: $C_k(p) = 0$.

Oltre al metodo per definire la credibilità di un link gli autori in [2] propongono un algoritmo, denominato *CredibleRank* di ranking basato sulla credibilità. *CredibleRank* definisce che la qualità di una pagina è determinata da due criteri: la qualità delle pagine che puntano ad essa e la credibilità di ogni pagina puntata. Un link da un alta-qualità/alta-credibilità conta più di un link da alta-qualità/bassa-credibilità. Definendo con $In(p)$ l'insieme di pagine che puntano a p . Calcoliamo *CredibleRank* $r_c(p)$ per una pagina p

$$r_c(p) = \sum_{q \in In(p)} C(q) \cdot r_c(q) \cdot w(q, p) \quad (2.13)$$

Questa formula dice che il valore di *CredibleRank* di una pagina p è determinato dalla qualità $r_c(q)$ e dalla credibilità dei link $C(q)$ delle pagine che la puntano così come la forza del link $w(q, p)$.

2.2.2 Metodi per identificare spam farm

Per riconoscere una spam farm si parte dal presupposto che i nodi della spam farm avranno dei link uscenti verso delle pagine target t per aumentarne il rank. In [7] per identificare le spam farm viene introdotta una misura, denominata *spam mass*, dell'impatto dello spam (basandosi sulla struttura del grafo) sul rank di una pagina. Le pagine target delle spam farm allora riceveranno, oltre ad un alto valore di *pagerank*, un alto valore di *spam mass* mentre le pagine non spam anche se hanno un alto valore di *pagerank* riceveranno un basso valore di *spam mass*. Un modo per stimare la *spam mass* per ogni nodo del grafo è partire dal presupposto che il web può essere partizionato in nodi non spam V^+ e nodi spam V^- e la loro unione forma il grafo del web. Per una data partizione $\{V^+, V^-\}$ di V e per dei nodi x il *pagerank* di x è la somma dei contributi di nodi non spam e dei nodi spam. Quindi vengono definite due misure di *spam mass*:

- La *spam mass assoluta* di x , denotata con M_x , è il *pagerank* che x riceve dai nodi spam è che uguale a:

$$M_x = q_x^{V^-} \quad (2.14)$$

dove $q_x^{V^-}$ è appunto il *pagerank* di x derivato dai nodi spam.

- La *spam mass relativa* di x , denotata da m_x , è la frazione del *pagerank* di x dovuto dal contributo dei nodi di spam cioè:

$$m_x = q_x^{V^-} / p_x \quad (2.15)$$

dove $q_x^{V^-}$ è il *pagerank* di x derivato dai nodi spam e p_x il *pagerank* derivato da tutti i nodi.

Dal momento che non è possibile conoscere le proprietà (spam o non spam) per tutti i nodi del grafo ma solo un sottoinsieme di nodi buoni $(\tilde{V})^+$ le misure precedenti vengono calcolate nel seguente modo:

- la stima assoluta di *spam mass* di un nodo x è:

$$\tilde{M}_x = p_x - p'_x \quad (2.16)$$

- la stima relativa di *spam mass* di x è:

$$\tilde{m}_x = (p_x - p'_x) / p_x = 1 - p'_x / p_x \quad (2.17)$$

dove $p = PR(v)$ è il *pagerank* dei nodi basato su una distribuzione uniforme mentre $p' = PR(v^{\tilde{V}^+})$ è *pagerank* basato sull'insieme $(\tilde{V})^+$ con una distribuzione di salto $v^{\tilde{V}^+}$. Nel caso in cui si conoscesse \tilde{V}^- lo *spam mass* può essere stimato con $M = PR(v^{\tilde{V}^-})$. Mentre se si conosceranno entrambi i sottoinsiemi V^+, V^- la stima dello *spam mass* può essere fatta attraverso $(\tilde{M} + M)/2$. Perciò è possibile utilizzare un valore di soglia tramite la quale una pagina è considerata facente parte di una spam farm se il valore di *spam mass* supera la soglia.

Bibliografia

- [1] Jacob Abernethy, Olivier Chapelle, and Carlos Castillo. Web spam identification through content and hyperlinks. In *Proceedings of the 4th International Workshop on Adversarial Information Retrieval on the Web*, AIRWeb '08, pages 41–44, New York, NY, USA, 2008. ACM.
- [2] James Caverlee and Ling Liu. Countering web spam with credibility-based link analysis. In *Proceedings of the Twenty-sixth Annual ACM Symposium on Principles of Distributed Computing*, PODC '07, pages 157–166, New York, NY, USA, 2007. ACM.
- [3] Cailing Dong and Bin Zhou. Effectively detecting content spam on the web using topical diversity measures. In *Proceedings of the The 2012 IEEE/WI-C/ACM International Joint Conferences on Web Intelligence and Intelligent Agent Technology - Volume 01*, WI-IAT '12, pages 266–273, Washington, DC, USA, 2012. IEEE Computer Society.
- [4] Nadav Eiron, Kevin S. McCurley, and John A. Tomlin. Ranking the web frontier. In *Proceedings of the 13th International Conference on World Wide Web*, WWW '04, pages 309–318, New York, NY, USA, 2004. ACM.
- [5] Dennis Fetterly, Mark Manasse, and Marc Najork. Spam, damn spam, and statistics: Using statistical analysis to locate spam web pages. In *Proceedings of the 7th International Workshop on the Web and Databases: Colocated with ACM SIGMOD/PODS 2004*, WebDB '04, pages 1–6, New York, NY, USA, 2004. ACM.

- [6] Dennis Fetterly, Mark Manasse, and Marc Najork. Detecting phrase-level duplication on the world wide web. In *Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '05, pages 170–177, New York, NY, USA, 2005. ACM.
- [7] Zoltan Gyongyi, Pavel Berkhin, Hector Garcia-Molina, and Jan Pedersen. Link spam detection based on mass estimation. In *Proceedings of the 32Nd International Conference on Very Large Data Bases*, VLDB '06, pages 439–450. VLDB Endowment, 2006.
- [8] Zoltan Gyongyi and Hector Garcia-Molina. Web spam taxonomy. Technical Report 2004-25, Stanford InfoLab, March 2004.
- [9] Zoltán Gyöngyi, Hector Garcia-Molina, and Jan Pedersen. Combating web spam with trustrank. In *Proceedings of the Thirtieth International Conference on Very Large Data Bases - Volume 30*, VLDB '04, pages 576–587. VLDB Endowment, 2004.
- [10] Nicholas R. Jennings. The global economic impact of spam. *Ferris Research*, 2005.
- [11] Nicholas R. Jennings. Cost of spam is flattening - our 2009 predictions. *Ferris Research*, 2009.
- [12] Jon M. Kleinberg. Authoritative sources in a hyperlinked environment. *J. ACM*, 46(5):604–632, September 1999.
- [13] Vijay Krishnan. Web spam detection with anti-trust rank. In *In AIRWEB*, pages 37–40, 2006.
- [14] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA, 2008. Pages 117–119.

- [15] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA, 2008. Pages 474–476.
- [16] Juan Martinez-Romo and Lourdes Araujo. Web spam identification through language model analysis. In *Proceedings of the 5th International Workshop on Adversarial Information Retrieval on the Web*, AIRWeb '09, pages 21–28, New York, NY, USA, 2009. ACM.
- [17] Alexandros Ntoulas, Marc Najork, Mark Manasse, and Dennis Fetterly. Detecting spam web pages through content analysis. In *Proceedings of the 15th International Conference on World Wide Web*, WWW '06, pages 83–92, New York, NY, USA, 2006. ACM.
- [18] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical Report 1999-66, Stanford InfoLab, November 1999. Previous number = SIDL-WP-1999-0120.
- [19] Stephen Robertson and Hugo Zaragoza. The probabilistic relevance framework: Bm25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389, April 2009.
- [20] Nikita Spirin and Jiawei Han. Survey on web spam detection: Principles and algorithms. *SIGKDD Explor. Newsl.*, 13(2):50–64, May 2012.
- [21] Tanguy Urvoy, Thomas Lavergne, and Pascal Filoche. Tracking web spam with hidden style similarity. In *AIRWeb*, pages 25–31, 2006.
- [22] Baoning Wu, Vinay Goel, and Brian D. Davison. Topical trustrank: Using topicality to combat web spam. In *Proceedings of the 15th International Conference on World Wide Web*, WWW '06, pages 63–72, New York, NY, USA, 2006. ACM.