

Spam-szűrés

Zongor Attila

Bevezetés

- Mi a spam?
- Miért kell szűrni a spameket?
- Hogyan védekezhetünk a spamek ellen?
- Hogyan működnek a spam szűrők?

Mi s spam?

- A spam legáltalánosabban kéretlen és tömeges e-mail üzenetet jelent.
- A címzett személye érdektelen, mivel az üzenetet változatlan formában lehetne sok más felhasználónak is küldeni.
- A címzett nem engedélyezte a levél küldését.
- A levél küldése és fogadása lényegesen nagyobb előnyt jelent a küldő számára

Miért kell szűrni a spamet?

- Nem szükség, viszont rengeteg időt takaríthatunk meg.
- Nem kell minden egyes levelet, mivel manapság naponta több száz spam is érkezhethet postaládánkba.

Hogyan védekezhetünk a spamek ellen?

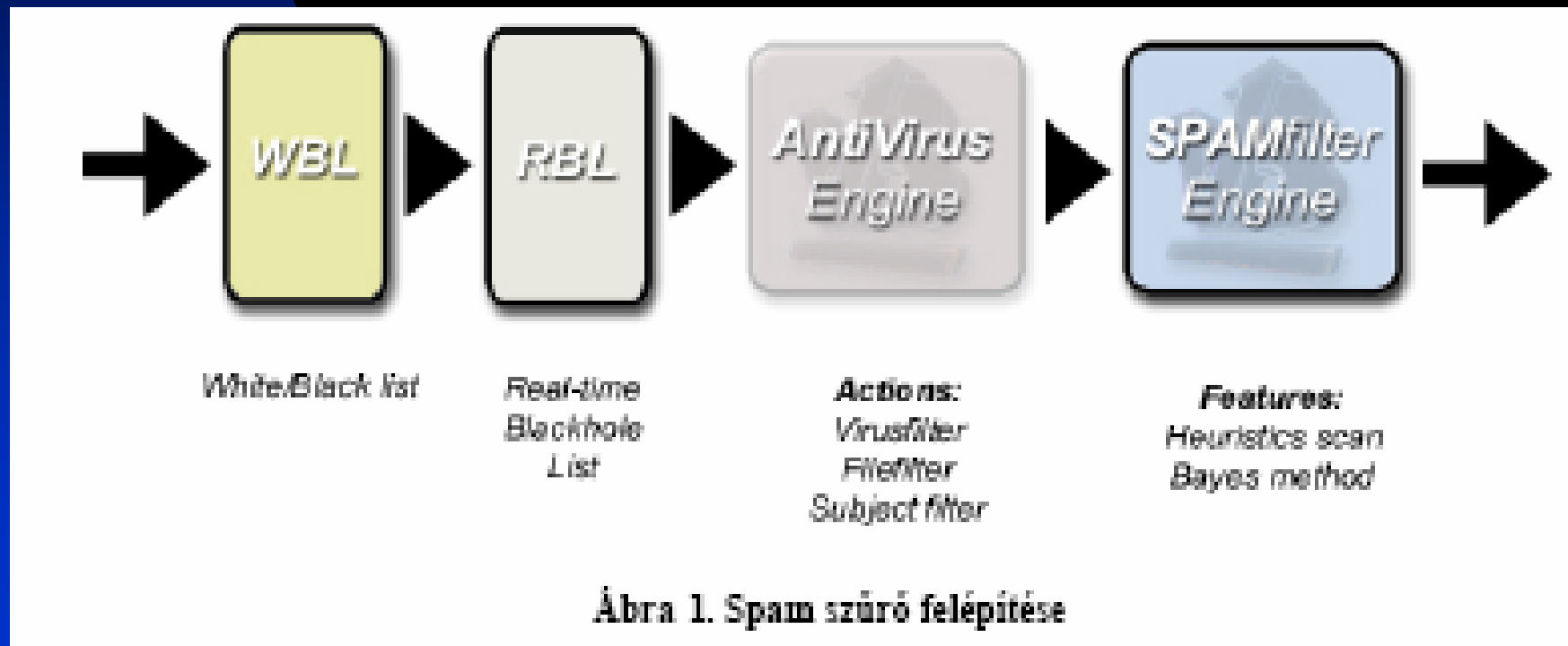
- Ne adjuk meg e-mail címünket ismeretlen helyen.
- Ne használjuk kedvenc email címünket elektronikus listákon.
- Nem nézem meg “EGYSZER AZ ÉLETBEN VAN ILYEN AKCIO” leveleket, ugyanis a levél megnyitásával a küldő megkapja az e-mail címünket amit továbbad...
- Használjuk a szolgáltató spam-szűrő rendszerét.

Hogyan működnek a spam szűrők?

- Statisztikai alapon:
 - Elemzik a levélek szövegét
 - Elemzik a levél fejlécét
 - Pár dolog ami spamre utal:
 - Nincs neve a feladónak, csak email címe.
 - Nem vagyunk közvetlen címzettek, azaz más is megkapta a levelet.
 - Nem normális SMTP szerverrel küldték a levelet, hanem helyi gépről.

Hogyan működnek a spam szűrők?

- Egy tipikus szűrő az alábbi séma szerint épül fel:



Hogyan működnek a spam szűrők?

- Az ellenőrzés első szintje: a white/black list
 - ★ A feladó címe vagy domain címe alapján történik a kiértékelés.
 - Ha white szerepel a cím a levél minden további ellenőrzés nélkül továbbításra kerül.
 - Ha black listán szerepel, akkor elutasításra kerül.
- A következő lépésben a spam szűrő megvizsgálja, hogy a küldő IP címe szerepel-e valamelyik folyamatosan frissített tiltólistán(RBL).
 - ★ Ezek a tiltólisták az ismert spamet küldő IP címeket és szolgáltatókat tartalmazzák.

Hogyan működnek a spam szűrők?

- A harmadik lépésben már a levél tartalmát elemzi a szűrő.
 - ★ Leggyakrabban használt módszer: BAYES-szűrésen alapuló döntéshozatal (annak a meghatározása, hogy egy adott levél milyen valószínűséggel spam). A feltételes valószínűségre vonatkozó Bayes tétel alapján vissza lehet vezetni annak a vizsgálatára, hogy a levél szavai milyen gyakorisággal fordulnak elő spam illetve nem spam levelekben.

Bayes szűrők

- Alapprobléma meghatározása, hogy az adott levél mekkora valószínűséggel spam.
- Bayes és a teljes valószínűség tétel alapján meghatározhatjuk.
- Sok levél feldolgozása után meghatározható hogy a szavak mekkora relatív gyakorisággal fordulnak elő spam illetve nem spam levelekben.
- Nagy mintalevél gyűjteményre van szükség...

Bayes szűrők

- Annak a valószínűsége, hogy a viagra szót tartalmazó levél spam (**spamicity**)
- Ahol $r(\text{viagra}|\text{spam})$ a viagra szónak spam levelekben való előfordulási aránya, $(\text{viagra}|\text{nospam})$ ugyan ez tiszta levelekben.
- $R(\text{spam})$ a spam minták aránya a tanító mintákban, értelemszerűen $r(\text{nospam}) = 1 - P(\text{spam}|\text{viagra})$

$$P(\text{spam} | \text{viagra}) \approx \frac{r(\text{viagra} | \text{spam}) \cdot r(\text{spam})}{r(\text{viagra} | \text{spam}) \cdot r(\text{spam}) + r(\text{viagra} | \text{nospam}) \cdot r(\text{nospam})}$$

Bayes szűrők

- Egy levél minden szavának meghatározza a **spamicity**-t.
- Kiválasztja a döntő szavakat és ezekre elvégzi az alábbi összesítést:

$$\text{Spamicity} = \frac{\prod_i \text{Spamicity}_i}{\prod_i \text{Spamicity}_i + \prod_i (1 - \text{Spamicity}_i)}$$

- Ha 1-hez közelít spam ha 0-hoz nem spam

Felismerés hatásfokának javítása

- Nem szavakat hanem szó-párokat elemzünk.
 - ◆Probléma: nagy adatbázis!
 - ◆Tanító adatbázis nagysága!

ROBS és ROBEX módszer

$$P(\text{spam} | \text{token}) = \frac{\text{robs} * \text{robx} + P(\text{token} | \text{spam})}{\text{robs} + P(\text{token} | \text{spam}) + P(\text{token} | \text{notspam}) * \text{scalefactor}}$$

- **Scalefactor** – spam és nem spam levelek hányadosa.
- **Robs** – mennyire számítson a robx és a token valószínűség.
- **Robx** – mennyi annak a val., hogy egy olyan token ami nem szerepel az adatbázisunkban, spam levélben szerepel.

ROBS és ROBEX módszer

$$P(sum) = \sum_{tokens} \left(\frac{P(token | spam)}{P(token | spam) + P(token | notspam) * scalefactor} \right)$$

- Megadja a spam valószínűségek összegét.

$$robex = \frac{P(sum)}{Count(spam) + Count(notspam)}$$

- Count(spam) –olyan tokenek száma amelyek csak spam levelekben szerepelnek.
- Count(notspam) –olyan tokenek amelyek csak tiszta levelekben fordulnak elő.

Vakriasztások kiküszöbölésének lehetőségei

- Hibák:
 - ★ Spam beengedések.
 - ★ Nem spam kizárása-
- Megtalálni a megfelelő arányt.

Vakriasztások kiküszöbölésének lehetőségei

- A spamnek minősítés határértékének növelése.
- Levél tanítás.
- Nem látott elemek valószínűségének beállítása.
- Minimum és maximum valószínűség.
 - ★ Megadjuk token minimum és maximum valószínűségét.

Felismerési és vakriasztási arány

	DNSBL	Mintakeresés	Heurisztika	Statisztika
Felismerés	0-60%	80%	95%	99%+
Vakriasztás	10%	2%	0.5%	0.1%

- DNSBL – nyilvántartott domain nevekre vonatkozó tiltólista.
- A statisztika alapon működő spam szűrők különböző felhasználók esetén felismerési arány csak:80-90%

Új spam esetén

- DNSBL – mivel a levél tartalma nem befolyásolja.....
- Mintakeresés – program illetve az adatbázis frissítésével tud lépéstartami.
- Heurisztika – hasonló, mint az előző.
- Statisztika – automatikusan képes megtanulni a spam jellemzőit abban az esetben, ha a felhasználó megadja, hogy melyik levél spam és melyik nem spam

Jellemzők szűrése

- Általános jellemzők: kis nagybetűs arány, betű és más karakter aránya.
- Spamekben használt trükkök felfedése: színek használata, szöveg elrejtése, véletlenszerű szövegek elhelyezése, stb..
- Levél feladójának ellenőrzése.
- Levél címzettjének ellenőrzése.
- Levél tárgyának vizsgálata.
- Stb.

Spamekben használt trükkök

- Képekben tárolt szöveg.
- Láthatatlan tinta
- Szavak tördelése több html tag-gel.
- Szöveg elrejtése HTML elé.
- Táblázatok.
- Szavak betűkre tördelése.

Program

- Megvalósítás
- Mit kódoltál,
- Milyen adatokra,
- Milyen eredménnyel