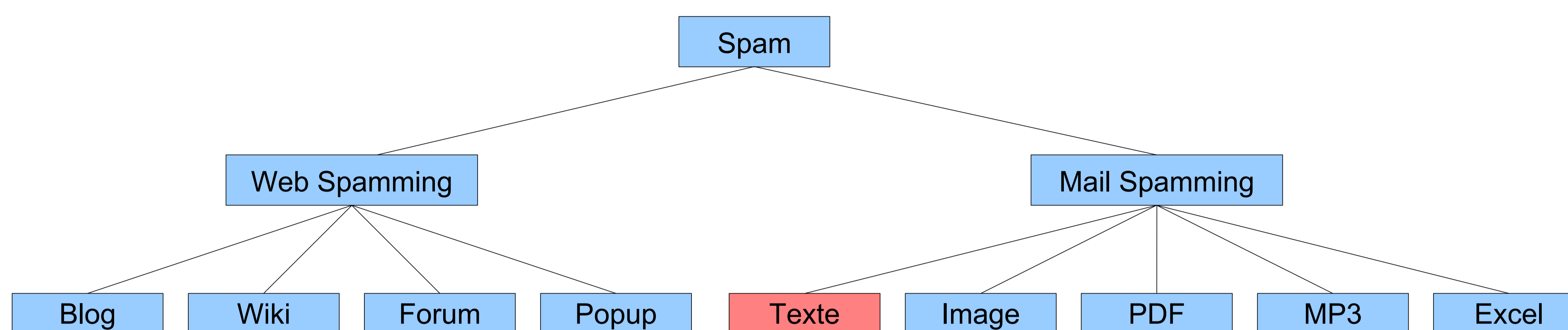


Les spams : la pollution informatique

Xavier D'Hondt, Vincent Stradiot, Jérémie Vion
Département d'Informatique

Le spam, ou pourriel est une communication électronique (mail, message, page internet...) non voulue par le receveur et envoyée en masse par l'expéditeur, en général à des fins publicitaires. La majorité de ces messages sont envoyés par courrier électronique. On estime que les spams représentent entre 90 et 97 % du volume total d'e-mails qui transitent à travers le web. Le terme spam est opposé au terme ham qui représente les « bons » mails.



Les différents types de spam

Un spammeur peut gagner entre 5000€ et 100 000€ par jour. Attention au revers de la médaille si le spammeur se fait attraper l'amende sera salée.



Certaines études démontrent que la pollution en dioxyde de carbone engendrée par les spams est équivalente à celle de 3 millions de voitures.



À ne pas faire :

- Répondre aux spams
- Transférer une chaîne à tous ses contacts
- Laisser son adresse à la vue de tout le monde
- Utiliser son adresse principale pour s'inscrire à des forums/concours

Filtre bayésien

Phase 1 : apprentissage

Phase 2 : filtrage

Exemple de filtrage

Supposons que durant la phase d'apprentissage, nous ayons rencontré 18 fois le mot « viagra » dans des spams et 2 fois dans des hams et que 80 % des messages reçus par l'utilisateur du filtre, jusqu'à présent, sont des spams.

a) Calculer la probabilité que le message soit un spam s'il contient « viagra »

$$P(\text{spam}|\text{viagra}) = \frac{18 \times 0,8}{18 \times 0,8 + 2 \times 0,2} = 0,97$$

b) Combiner les résultats afin d'obtenir un résultat pour le message entier (exemple : « viagra (0,97) améliore (0,41) performances (0,62) sexuelles (0,91) »)

$$P(\text{spam}) = \frac{0,97 \times 0,41 \times 0,62 \times 0,91}{0,97 \times 0,41 \times 0,62 \times 0,91 + 0,03 \times 0,59 \times 0,38 \times 0,09} \approx 0,99$$

c) Comparer la valeur obtenue à la valeur seuil prédéfinie

$$0,99 \geq 0,95$$

d) Classifier le message comme spam ou non

