

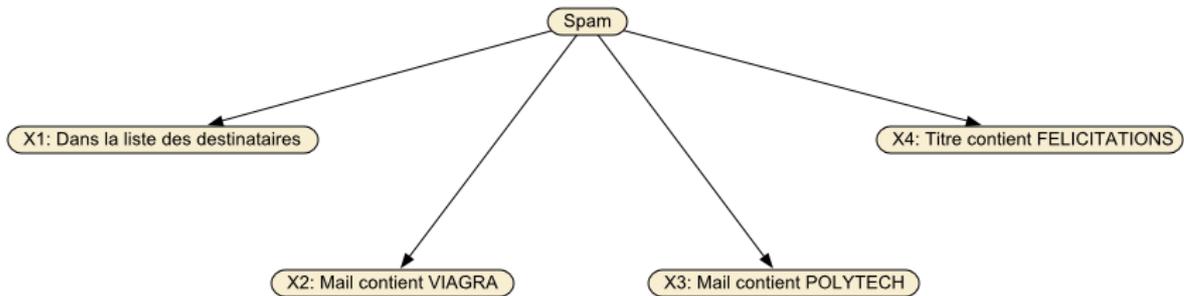
Systèmes de Raisonnement Probabiliste – Examen

Durée : 1h30

Documents non autorisés - Calculatrice autorisée

Bayesian spam filtering

De nombreux anti-spams utilisent des modèles basés sur un réseau bayésien très simple, le classifieur de Bayes naïf, dont la structure est donnée ci-dessous.



Les différentes valeurs prises par les variables du modèle sont : Spam = oui / non, X1 = seul/ parmi N / non, X2 = oui / non, X3 = oui / non, X4 = oui / non. Une analyse utilisant une base conséquente de mails (100 spams et 400 légitimes) a permis de déterminer que :

- Vous êtes très rarement seul destinataire d'un spam (2%). Le plus souvent vous n'apparaissez pas dans la liste des destinataires (78%). Ces proportions sont inversées lorsque le mail n'est pas un spam.
- 1% des spams contiennent le mot POLYTECH, et 10% le mot VIAGRA. 30% des spams sont reconnaissables par un titre contenant le mot FELICITATIONS. A l'inverse, 30% des mails légitimes contiennent le mot POLYTECH, 0.1% le mot VIAGRA, et 1% ont un titre contenant le mot FELICITATIONS.

1. (2 points) Propriétés du modèle

- (a) (1 point) Quelles sont les propriétés de dépendance et d'indépendance entre X1 et Spam ? entre X1 et X2 ? entre X1 et X2 conditionnellement à Spam ?
- (b) (1 point) Quelle est la décomposition de la loi jointe pour ce modèle.

2. (8 points) Message Passing

- (a) (3 points) Quelle est la probabilité qu'un mail soit un Spam sachant que vous êtes le seul destinataire ? si vous faites partie de la liste des destinataires ? et si vous n'apparaissez pas dans la liste ?
- (b) (3 points) Quelle est la probabilité qu'un mail soit un Spam sachant que vous êtes un des destinataires (mais pas le seul) et que le texte contient le mot VIAGRA ?
- (c) (2 points) Même question en rajoutant le fait que le mot POLYTECH n'apparaît pas et que le titre contient FELICITATIONS.

3. (5 points) Optimisation

Le classifieur de Bayes naïf possède une propriété intéressante permettant d'accélérer les calculs d'inférence.

- (a) (3 points) Montrez que $P(\text{Spam}|X_1X_2X_3X_4) \propto \prod_i P(\text{Spam}|X_i)$
- (b) (2 points) Vérifiez-le par le calcul en reprenant l'exemple du mail précédent (2.c).

4. (5 points) Naïf augmenté

- (a) (3 points) Supposons maintenant qu'en réalité X3 et X4 sont dépendants conditionnellement à Spam. Proposer un nouveau modèle (structure et tables de probabilités) prenant cette information en compte. Quelle est la loi jointe de ce nouveau modèle
- (b) (2 points) Proposez un ou plusieurs algorithmes d'inférence permettant d'utiliser ce modèle. Justifiez votre réponse.