

# Utilisation de la Machine Cellulaire pour la Détection des Courriels Indésirables

F. Barigou\*, B. Atmani\*\*, B. Beldjilali\*\*\*

Equipe SIF, Laboratoire d'Informatique d'Oran

Université d'Oran BP 1524, El M'Naouer, 31 000 Oran, Algérie

\* fatbarigou@gmail.com, \*\*atmani.baghdad@univ-oran.dz, \*\*\*bouzianebedjilali@yahoo.fr

## 1 Introduction

Dans ce papier, nous proposons pour la première fois une approche de filtrage de spam qui se base sur l'induction symbolique par automate cellulaire (Atmani et Beldjilali, 2007). Ce choix de cette technique a été motivé par ses propriétés intéressantes comme la réduction de l'espace de stockage et le temps de classification. Le principe de cette approche est très simple, il s'agit de construire un modèle booléen à partir d'un ensemble de courriels d'apprentissage. Ce modèle, qui sera utilisé pendant la phase de classification, va permettre de déterminer la nature d'un nouveau courriel (spam ou légitime).

## 2 Approche Cellulaire de Classification

L'automate cellulaire CASI (Cellular Automata for System Induction) issue des travaux de (Atmani et Beldjilali, 2007) est une méthode cellulaire de génération, de représentation et d'optimisation des graphes d'induction (Zighed,2000) générés à partir d'un ensemble d'exemples d'apprentissage. Ce système cellulo-symbolique est organisé en cellules où chacune d'elles, est reliée seulement avec son voisinage. Toutes les cellules obéissent en parallèle à la même règle appelée fonction de transition locale, qui a comme conséquence une transformation globale du système. Deux composants coopèrent entre eux pour la construction du modèle booléen : le COG (Cellular Optimization and Generation) qui s'occupe de la génération du graphe d'induction cellulaire et de son optimisation et le CIE (Cellular Inference Engine), un moteur d'inférence cellulaire, qui génère un ensemble de règles cellulaires sous formes conjonctives utilisées pendant la phase de filtrage. Pour se faire, ils utilisent une base de connaissances sous forme de deux couches finies d'automates finis. La première couche, CelFact<sup>1</sup>, pour la base des faits et, la deuxième couche, CelRule<sup>2</sup>, pour la base de règles. Le voisinage des cellules est défini par deux matrices d'incidence d'entrée  $R_E$  et de sortie  $R_S$ . La dynamique de l'automate cellulaire, utilise deux fonctions de transitions  $\delta_{fact}$  qui simule les phases de sélection et de filtrage dans un système expert et  $\delta_{rule}$  qui correspond à la phase d'exécution :

$$\begin{aligned} (EF, IF, SF, ER, IR, SR) &\xrightarrow{\delta_{fact}} (EF, IF, EF, ER + (R_E^T * EF), IR, SR) \\ (EF, IF, SF, ER, IR, SR) &\xrightarrow{\delta_{rule}} (EF + (R_S * ER), IF, SF, ER, IR, \overline{ER}) \end{aligned}$$

<sup>1</sup> Toute cellule de CelFact est considérée comme fait établi si sa valeur est 1, sinon, elle est considérée comme fait à établir. Elle se présente sous trois états : état d'entrée (EF), état interne (IF) et état de sortie (SF)

<sup>2</sup> Toute cellule de CelRule est considérée comme une règle candidate si sa valeur est 1, sinon, elle est considérée comme une règle qui ne doit pas participer à l'inférence. Elle se présente sous trois états : état d'entrée (ER), état interne (IR) et état de sortie (SR)