

# Informatique Décisionnelle (5073)

## Session de Septembre 2018

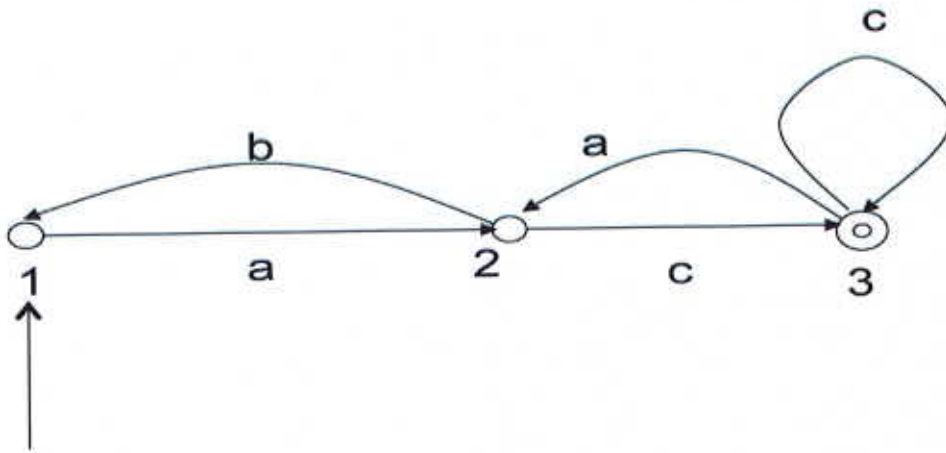
4 pages

Aucun document n'est autorisé. Les cinq problèmes sont indépendants.

### 1. Problème de mots.

Soit l'automate  $A$  défini par  $Q=\{1,2,3\}$  où 1 est l'état initial et 3 est l'état final, comme indiqué sur la figure ci-dessous. Les transitions sont définies par le graphe de la figure.

- Trouver un mot de longueur au moins 6, accepté et un autre rejeté.
- Décrire le langage défini par cet automate par une expression régulière.
- Préciser si l'automate est déterministe ou non déterministe.
- Pour la distance d'édition, quelle est la distance du mot  $bbbcccbcb$  au langage de l'automate.
- Trouver un automate  $B$ , dont le langage soit  $\varepsilon$ -loin de  $A$ .
- Etant donné un mot  $w$ , expliquer comment le classifier comme un mot de  $A$  ou comme un mot de  $B$ , selon la distance d'édition.



## 2. Problème d'arbres.

Soit le fichier XML suivant :

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<MEDIA>
  <FILM LANG="fr">
    <AUTHOR>
      <FIRSTNAME>Alain</FIRSTNAME>
      <LASTNAME>Michard</LASTNAME>
    </AUTHOR>
    <TITLE>Le tigre</TITLE>
    <ACTOR>
      <FIRSTNAME>Alain</FIRSTNAME>
      <LASTNAME>Delon</LASTNAME>
    </ACTOR>
    <ACTOR>
      <FIRSTNAME>Catherine</FIRSTNAME>
      <LASTNAME>Deneuve</LASTNAME>
    </ACTOR>
    <DATEPUB>1998</DATEPUB>
  </FILM>
  <FILM LANG="an">
    <AUTHOR>
      <FIRSTNAME>William J.</FIRSTNAME>
      <LASTNAME>Pardi</LASTNAME>
    </AUTHOR>
    <TITLE>Le lion</TITLE>
    <PUBLISHER>
      <NAME>MGM</NAME>
      <PLACE>Hollywood</PLACE>
    </PUBLISHER>
    <DATEPUB>1999</DATEPUB>
  </FILM>
</MEDIA>
```

- Décrire ce fichier par un arbre XML, en abrégiant le nom des tags, par leur première lettre, sans écrire les attributs. Utiliser les deux premières lettres en cas d'ambiguïté.  
<AUTHOR> sera considéré comme la balise <au> et <ACTOR> sera considéré comme la balise <ac>.
- Trouver une DTD pour laquelle cet arbre est valide.
- Quel est le langage d'arbres défini par cette DTD.
- Trouver un arbre proche, mais non valide, pour cette DTD.
- Trouver une DTD<sub>1</sub>, qui soit ε-loin de cette DTD pour la distance d'édition.

### 3. Problème de Schéma relationnel.

Le système d'information du métro est une base de données relationnelles, qui se résume à 4 tables. Les tables **TRAINS**, **STATIONS**, **RESEAU** et **TRAFIC**.

**TRAINS**(NumT, Marque, Capacité, Dépôt, Année) décrit les trains avec un identifiant (NumT), une Marque (le constructeur), un nombre de passagers (Capacité), un Dépôt (Lieu), et l'année de mise ne service du véhicule (Année).

**STATIONS**(NumS, Nom, Adresse, Arrondissement) décrit une station de métro avec un identifiant (NumS), un Nom, une Adresse, un Arrondissement.

**RESEAU** (Départ, Arrivée, Ligne) décrit les segments de chaque ligne de métro. Si une ligne décrit un chemin sur 20 stations, il y a 19 segments et donc 19 tuples. Le tuple (Odéon, St-Michel, 4) décrit un segment de la ligne 4

#### TRAINS

<u>NumT</u>	Marque	Capacité	Dépôt	Année
-------------	--------	----------	-------	-------

#### STATIONS

<u>NumS</u>	Nom	Adresse	Arrondissement
-------------	-----	---------	----------------

#### RESEAU

Départ	Arrivée	Ligne
--------	---------	-------

La table **TRAFIC** décrit les différents trajets effectués par les trains. Les attributs de la table **TRAFIC** décrivent le trajet d'un train (*NumT*), qui part à une *Heure Départ*, d'une *Station Départ*, à une certaine *Date*, pour arriver à un *Lieu Arrivée*, à une *Heure d'Arrivée Prévue*, sur une *Ligne* du métro, avec peut-être un *Retard*, et/ou un *Incident* (0 s'il n'y a pas s'incident et 1 s'il y a un incident).

#### TRAFIC

NumT	Heure Départ	Station Départ	Date	Lieu Arrivée	Heure Arrivée Prévue	Ligne	Retard	Incident
------	-----------------	-------------------	------	-----------------	----------------------------	-------	--------	----------

- a. Décrire un schéma Entité-relations associé à ce Schéma relationnel.
- b. Ecrire des requêtes SQL pour trouver :
  1. Dans quel arrondissement se trouve la station « Chatelet » ?
  2. Quelle est la marque des trains qui ont circulé à « Chatelet » le 1<sup>er</sup> juin 2018 ?
  3. Combien de lignes (distinctes) desservent la station « Chatelet » ?

4. Combien de train ont circulé le 20 septembre 2017 sur la ligne 4 ?
- c. Existe-t-il une requête SQL unaire (avec seul un attribut NumS) qui donnerait toutes les stations accessibles depuis la station « Chatelet » sur la ligne 4 ?
- d. Indiquer les dépendances fonctionnelles des tables **TRAFIC** et **RESEAU**.
- e. Quelle est la clé de la table **TRAFIC** ?

#### 4. Problème OLAP.

On souhaite analyser la table **TRAFIC** du schéma relationnel précédent, considérée comme un entrepôt de données pour comprendre la répartition du trafic, des retards et des incidents.

- a. Proposez un schéma étoile pour la table **TRAFIC**, avec *Retards* et *Incidents* comme mesures (2 mesures) et Somme comme Agrégation.
- b. Donnez les expressions OLAP (chemin, mesure, agrégation) associées aux requêtes :
  - Analyse des Incidents par *Lignes* en 2006
  - Analyse des Incidents par *Marque* des Trains,
  - Analyse des Retards en Avril 2014, par *Ligne* et par *Jour*.
- c. Donnez une expression en SQL-étendu (utilisant GROUP BY) pour la 1<sup>ère</sup> requête.

#### 5. Fouille de Données

L'entrepôt de données (la table TRAFIC du problème 4) est alimenté par les données de la RATP, selon le schéma du problème 3. On cherche à prédire les incidents, à partir de N=1000 enregistrements de cette table.

- a. On discrétise la variable *Retard* à l'aide d'une nouvelle variable R :  $R=0$  si  $Retard < 5mins$ ,  $R=1$  si  $5mins < Retard < 15mins$ , et  $R=2$  si  $Retard > 15mins$ . Donner un arbre de décision pour prédire le Retard
- b. Expliquer comment construire un arbre de décision pour la table **TRAFIC**, en utilisant le critère du Gain d'Information. Rappeler la définition de l'erreur pour prédire s'il existe un incident (Incident=1) ou non (Incident=0). Donner un exemple d'arbre de décision dans ce contexte et son taux d'erreur.
- c. Même question que a, à l'aide d'une régression logistique. Expliquer comment construire une régression logistique.