

---

---

# Données textuelles

---

---

*Dan VODISLAV*

**CY Cergy Paris Université**  
**Master Informatique M1**  
**Cours BDA**

---

---

## Plan

---

---

- Traitement des données textuelles
- Indexation du texte
- Requêtes mots-clés: modèles booléen et vecteur

## Données textuelles

---

- **Problème:** comment gérer les données de type texte, afin de répondre efficacement à des requêtes mots-clés?
- **Exemple (en anglais):** ensemble de 7 documents
  - d1:* The jaguar is a New World mammal of the Felidae family.
  - d2:* Jaguar has designed four new engines.
  - d3:* For Jaguar, Atari was keen to use a 68K family device.
  - d4:* The Jacksonville Jaguars are a professional US football team.
  - d5:* Mac OS X Jaguar is available at a price of US \$199 for Apple's new "family pack".
  - d6:* One such ruling family to incorporate the jaguar into their name is Jaguar Paw.
  - d7:* It is a big cat.

## Pré-traitement du texte

---

- **Étapes de pré-traitement** dépendantes de l'application et de la langue des documents
  - Découpage en mots ("tokenization")
  - Lemmatisation ("stemming")
  - Élimination des mots fréquents ("stop words")

# Découpage en mots

---

---

- Le mot est l'unité de recherche
- Le découpage en mots n'est pas si simple que ça!
  - Dans certaines langues, les mots ne sont pas séparés par des espaces (chinois, japonais)
  - Traitement spécial pour les acronymes, les abréviations, les élisions, les nombres, les URLs, les adresses mail, les unités de mesure, etc.
  - Mots composés (hostname, host-name, host name): les garder ensemble ou les séparer en mots?
    - Parfois très compliqué (Allemand): analyse lexicale et linguistique
- Dans cette étape: on enlève la ponctuation et on passe en minuscules

## Exemple découpage en mots

---

---

*d1*: the<sub>1</sub> jaguar<sub>2</sub> is<sub>3</sub> a<sub>4</sub> new<sub>5</sub> world<sub>6</sub> mammal<sub>7</sub> of<sub>8</sub> the<sub>9</sub> felidae<sub>10</sub> family<sub>11</sub>  
*d2*: jaguar<sub>1</sub> has<sub>2</sub> designed<sub>3</sub> four<sub>4</sub> new<sub>5</sub> engines<sub>6</sub>  
*d3*: for<sub>1</sub> jaguar<sub>2</sub> atari<sub>3</sub> was<sub>4</sub> keen<sub>5</sub> to<sub>6</sub> use<sub>7</sub> a<sub>8</sub> 68k<sub>9</sub> family<sub>10</sub> device<sub>11</sub>  
*d4*: the<sub>1</sub> jacksonville<sub>2</sub> jaguars<sub>3</sub> are<sub>4</sub> a<sub>5</sub> professional<sub>6</sub> us<sub>7</sub> football<sub>8</sub> team<sub>9</sub>  
*d5*: mac<sub>1</sub> os<sub>2</sub> x<sub>3</sub> jaguar<sub>4</sub> is<sub>5</sub> available<sub>6</sub> at<sub>7</sub> a<sub>8</sub> price<sub>9</sub> of<sub>10</sub> us<sub>11</sub> \$199<sub>12</sub> for<sub>13</sub>  
apple's<sub>14</sub> new<sub>15</sub> family<sub>16</sub> pack<sub>17</sub>  
*d6*: one<sub>1</sub> such<sub>2</sub> ruling<sub>3</sub> family<sub>4</sub> to<sub>5</sub> incorporate<sub>6</sub> the<sub>7</sub> jaguar<sub>8</sub> into<sub>9</sub> their<sub>10</sub>  
name<sub>11</sub> is<sub>12</sub> jaguar<sub>13</sub> paw<sub>14</sub>  
*d7*: it<sub>1</sub> is<sub>2</sub> a<sub>3</sub> big<sub>4</sub> cat<sub>5</sub>

# Lemmatisation

---

---

- Considérer la racine (lemme) des mots, afin de "fusionner" tous les mots de la même famille
  - Pas toutes les applications souhaitent cela!
  - En cherchant *enfant* on retrouve *enfance* ou *enfantin*
  - Différents degrés de lemmatisation
  - Il est possible de construire plusieurs index, avec des techniques de lemmatisation différentes
- Types de lemmatisation
  - Morphologique
  - Lexicale
  - Phonétique

## Types de lemmatisation

---

---

- Morphologique: enlever les terminaisons
  - Pluriel, genre, temps, mode, ...
  - Pas toujours facile: "*Les poules du couvent couvent.*"
  - En anglais c'est plus facile (terminaisons régulières, avec quelques exceptions bien répertoriées), mais il reste des ambiguïtés
- Lexicale
  - Termes d'une même famille
  - En anglais: algorithme de Porter, basé sur des critères morphologiques (!)
    - Pas toujours de bons résultats (ex. *university*, *universal* → *univers*)
  - Couplage avec des dictionnaires pour regrouper aussi les synonymes
- Phonétique
  - Regroupe des mots qui se prononcent pareil (ou presque pareil)
  - Objectif: traiter les fautes de frappe/orthographe
  - Assez grossier

## Exemple lemmatisation

---

*d1*: the<sub>1</sub> jaguar<sub>2</sub> be<sub>3</sub> a<sub>4</sub> new<sub>5</sub> world<sub>6</sub> mammal<sub>7</sub> of<sub>8</sub> the<sub>9</sub> felidae<sub>10</sub> family<sub>11</sub>

*d2*: jaguar<sub>1</sub> have<sub>2</sub> design<sub>3</sub> four<sub>4</sub> new<sub>5</sub> engine<sub>6</sub>

*d3*: for<sub>1</sub> jaguar<sub>2</sub> atari<sub>3</sub> be<sub>4</sub> keen<sub>5</sub> to<sub>6</sub> use<sub>7</sub> a<sub>8</sub> 68k<sub>9</sub> family<sub>10</sub> device<sub>11</sub>

*d4*: the<sub>1</sub> jacksonville<sub>2</sub> jaguar<sub>3</sub> be<sub>4</sub> a<sub>5</sub> professional<sub>6</sub> us<sub>7</sub> football<sub>8</sub> team<sub>9</sub>

*d5*: mac<sub>1</sub> os<sub>2</sub> x<sub>3</sub> jaguar<sub>4</sub> be<sub>5</sub> available<sub>6</sub> at<sub>7</sub> a<sub>8</sub> price<sub>9</sub> of<sub>10</sub> us<sub>11</sub> \$199<sub>12</sub> for<sub>13</sub>  
apple<sub>14</sub> new<sub>15</sub> family<sub>16</sub> pack<sub>17</sub>

*d6*: one<sub>1</sub> such<sub>2</sub> rule<sub>3</sub> family<sub>4</sub> to<sub>5</sub> incorporate<sub>6</sub> the<sub>7</sub> jaguar<sub>8</sub> into<sub>9</sub> their<sub>10</sub> name<sub>11</sub>  
be<sub>12</sub> jaguar<sub>13</sub> paw<sub>14</sub>

*d7*: it<sub>1</sub> be<sub>2</sub> a<sub>3</sub> big<sub>4</sub> cat<sub>5</sub>

---

## Élimination des mots fréquents

---

- Les mots fréquents se trouvent partout et ne distinguent pas les documents les uns par rapport aux autres
  - Leur élimination ne modifie pas significativement les résultats des requêtes
- Avantage: le nombre de mots du texte et la taille de l'index diminuent sensiblement
- Mots fréquents:
  - Articles: le, la, un, une, des, ...
  - Verbes fonctionnels: être, avoir, faire, ...
  - Conjonctions: et, ou, que, ...
  - etc.

## Exemple élimination mots fréquents

---

*d1*: jaguar<sub>2</sub> new<sub>5</sub> world<sub>6</sub> mammal<sub>7</sub> felidae<sub>10</sub> family<sub>11</sub>

*d2*: jaguar<sub>1</sub> design<sub>3</sub> four<sub>4</sub> new<sub>5</sub> engine<sub>6</sub>

*d3*: jaguar<sub>2</sub> atari<sub>3</sub> keen<sub>5</sub> 68k<sub>9</sub> family<sub>10</sub> device<sub>11</sub>

*d4*: jacksonville<sub>2</sub> jaguar<sub>3</sub> professional<sub>6</sub> us<sub>7</sub> football<sub>8</sub> team<sub>9</sub>

*d5*: mac<sub>1</sub> os<sub>2</sub> x<sub>3</sub> jaguar<sub>4</sub> available<sub>6</sub> price<sub>9</sub> of<sub>10</sub> us<sub>11</sub> \$199<sub>12</sub> apple<sub>14</sub> new<sub>15</sub>  
family<sub>16</sub> pack<sub>17</sub>

*d6*: one<sub>1</sub> such<sub>2</sub> rule<sub>3</sub> family<sub>4</sub> to<sub>5</sub> incorporate<sub>6</sub> jaguar<sub>8</sub> their<sub>10</sub> name<sub>11</sub> jaguar<sub>13</sub>  
paw<sub>14</sub>

*d7*: big<sub>4</sub> cat<sub>5</sub>

---

## Index inversé

---

- Index de tous les mots (termes) retenus après le pré-traitement, avec pour chacun la liste de documents où il apparaît
- Stockage
  - Petite échelle: sur disque, avec des mécanismes de mémoire virtuelle
  - Grande échelle: *cluster de machines*, avec hachage pour trouver la machine qui stocke l'entrée pour un mot donné
- Mise à jour
  - Coûteuse, réalisée généralement par lots (batch) et non pas individuellement

## Exemple index inversé

---

**family:** *d1, d3, d5, d6*  
**football:** *d4*  
**jaguar:** *d1, d2, d3, d4, d5, d6*  
**new:** *d1, d2, d5*  
**rule:** *d6*  
**us:** *d4, d5*  
**world:** *d1*  
...

## Positions des mots dans le document

---

- Utiles pour les requêtes de type "phrase" (plusieurs mots qui doivent se suivre dans l'ordre) ou des opérateurs comme NEAR
- Rajoutées à l'index pour chaque document
  - Un document peut avoir plusieurs occurrences du mot

**family:** *d1/11, d3/10, d5/16, d6/4*  
**football:** *d4/8*  
**jaguar:** *d1/2, d2/1, d3/2, d4/3, d5/4, d6/8+13*  
**new:** *d1/5, d2/5, d5/15*  
**rule:** *d6/3*  
**us:** *d4/7, d5/11*  
**world:** *d1/6*  
...

## Requêtes booléennes

---

- Trouver les documents qui contiennent les mots de la requête
- Un seul mot: liste de documents dans l'index inversé
- Plusieurs mots
  - Ex. (*jaguar AND new AND NOT family*) OR *cat*
  - On cherche les listes de documents pour les mots de la requête
  - Opérations sur les listes, suivant l'opérateur logique:
    - *AND*: intersection
    - *OR*: union
    - *AND NOT*: différence
- Requêtes impliquant la position dans le texte
  - Pareil, mais on filtre par rapport à la position des mots

## Modèle vecteur

---

- Réponse: documents *similaires* à une requête
- Méthode
  - Dans chaque document une partie des mots sont indexés
    - $n$ : nombre total de mots indexés dans le système
  - Un document  $d$  devient un point (vecteur) dans un espace à  $n$  dimensions
    - à chacun des mots ( $i$ ) est associé un poids (une coordonnée)  $1 \geq w_{id} \geq 0$
  - Requête  $q$  : point (vecteur):  $(w_{1q}, \dots, w_{iq}, \dots, w_{nq})$
  - Réponse: les documents (points) les plus proches du point requête
- Avantage: ordonnancement des réponses par similarité décroissante
- Modèle booléen: les poids sont 0 / 1 et l'appartenance d'un document au résultat est stricte

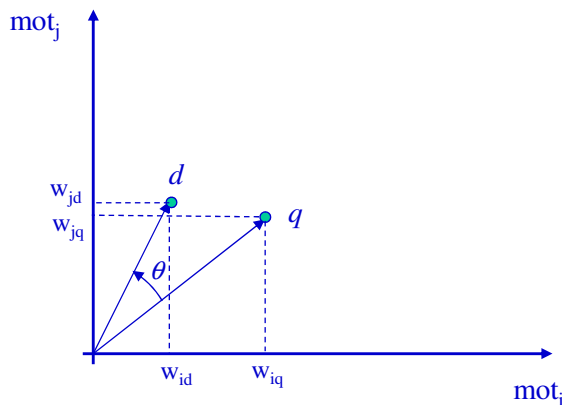


# Modèle vecteur : similarité

---

- Exemples de mesures de similarité

- Distance euclidienne
- On préfère: cosinus de l'angle entre les deux vecteurs
  - $\text{sim}(d, q) = \cos(\theta) = \sum (w_{id} * w_{iq}) / (\text{norme}(d) * \text{norme}(q))$



## Calcul des poids: tf/idf

---

- Fréquence d'un terme (mot) dans un document

- $tf$ : « term frequency »
- exprime l'importance du terme ( $i$ ) pour caractériser le document ( $d$ )

$$tf_{id} = \text{freq}_{id} / \max_t(\text{freq}_{td})$$

- Inverse de la fréquence d'un terme dans un corpus de documents

- $idf$ : « inverse document frequency »
- exprime la capacité du terme ( $i$ ) à filtrer parmi les documents du corpus
  - corpus de  $N$  documents, le terme  $i$  apparaît dans  $n_i$  documents

$$idf_i = \log(N/n_i)$$

- Calcul des poids :  $w_{id} = tf_{id} * idf_i$

- $w_{iq}$  : une formule similaire (avec des variantes), en considérant  $q$  comme un document contenant les mots de la requête

## Exemple calcul tf/idf

---

**family:** *d1/11, d3/10, d5/16, d6/4*  
**football:** *d4/8*  
**jaguar:** *d1/2, d2/1, d3/2, d4/3, d5/4, d6/8+13*  
**new:** *d1/5, d2/5, d5/15*  
**rule:** *d6/3*  
**us:** *d4/7, d5/11*  
**world:** *d1/6*

...

- $\mathbf{tf_{jaguar,d6} = freq_{jaguar,d6} / \max_t(freq_{t,d6}) = 2/2 = 1}$   
 $\mathbf{tf_{rule,d6} = freq_{rule,d6} / \max_t(freq_{t,d6}) = 1/2 = 0,5}$   
 $\mathbf{tf_{new,d6} = freq_{new,d6} / \max_t(freq_{t,d6}) = 0/2 = 0}$
- $\mathbf{idf_{jaguar} = \log(N/n_{jaguar}) = \log(7/6) = 0,067}$   
 $\mathbf{idf_{rule} = \log(N/n_{rule}) = \log(7/1) = 0,845}$   
 $\mathbf{idf_{new} = \log(N/n_{new}) = \log(7/3) = 0,368}$
- $\mathbf{W_{jaguar,d6} = tf_{jaguar,d6} * idf_{jaguar} = 0,067}$   
 $\mathbf{W_{rule,d6} = tf_{rule,d6} * idf_{rule} = 0,422}$   
 $\mathbf{W_{new,d6} = 0}$