

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

UVF3B403 MS IABDA

Fouille de texte

Yannis Haralambous (IMT Atlantique)

27 janvier 2017

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- La *fouille de texte* (en anglais *text mining*) est la découverte de connaissances nouvelles dans des données textuelles (potentiellement de volume important).
- On y combine des techniques de fouille de données, d'apprentissage artificiel, de statistique, de traitement automatique de langue.
- Avec la montée en puissance des machines toute méthode de traitement de langue peut potentiellement être utilisée pour la fouille (exemple : la traduction automatique → alignement de textes → fouille multilingue, désambiguïsation, etc.).

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

Nombreuses applications :

- sécurité : trouver des terroristes par leur tweets... ;
- veille concurrentielle : analyse des articles de presse, brèves, blogs, forums, Twitter, etc. (opinion) ;
- veille scientifique : analyse et résumé automatique des publications, des brevets, des pages Web, des blogs, etc. (découvertes) ;
- économie / marketing : gestion de relation clients, analyse de messages ;
- systèmes question-réponse en langue naturelle (achats en ligne, demandes de renseignements) ;
- extraction d'information (→ Web mining, etc.) ;
- classification de textes (par exemple : SPAM) ;
- découverte scientifique, en particulier dans le secteur biomédical (40 kpublications par mois) ;
- ... tout ce qui relève de l'utilisation de la langue sur le Web, et ailleurs.

La fouille de texte - le biomédical

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

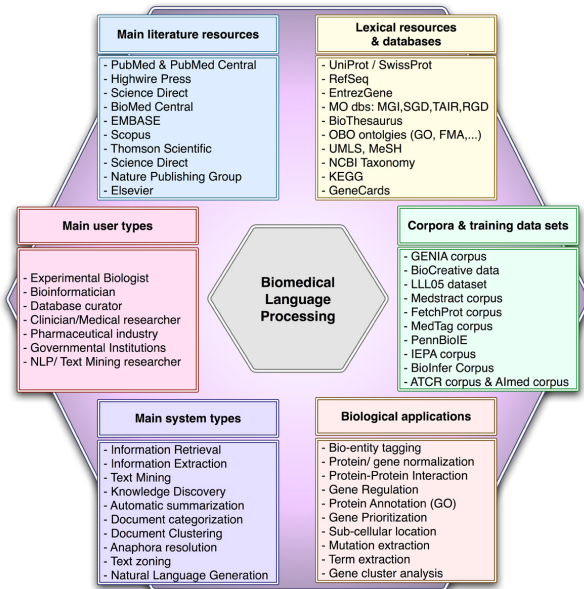
Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes





La fouille de texte - le biomédical

UVF3B403 MS

IABDA

Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

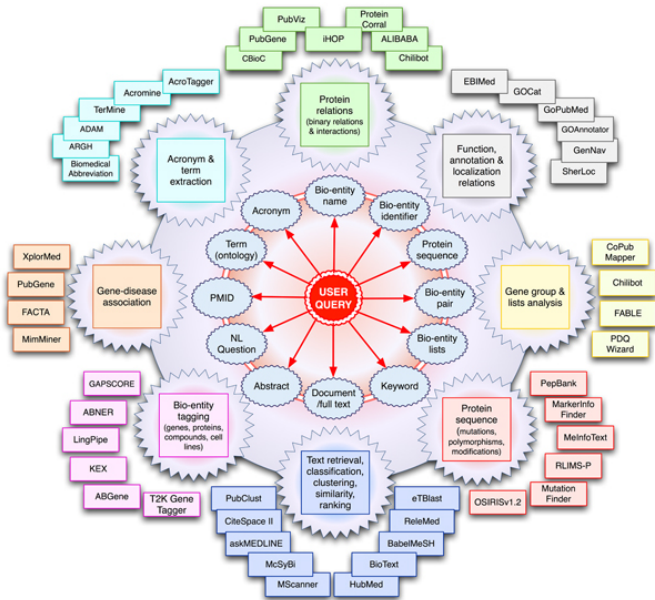
Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

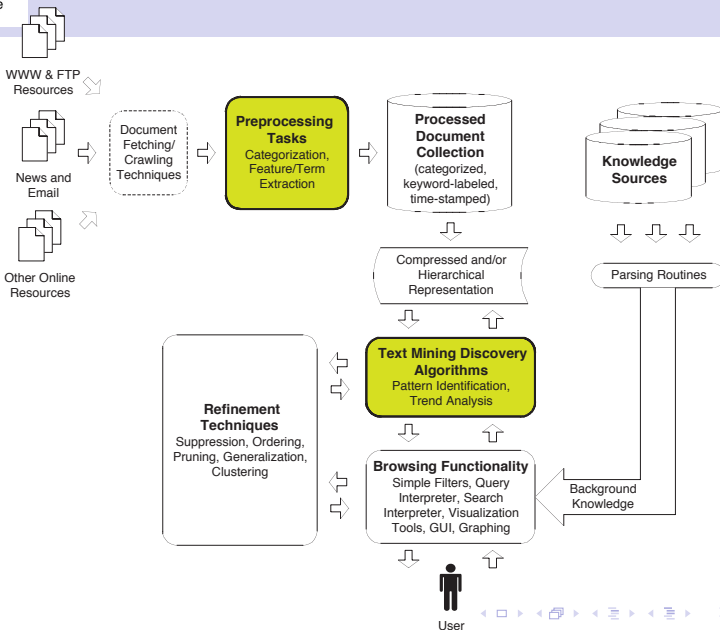
Requêtes



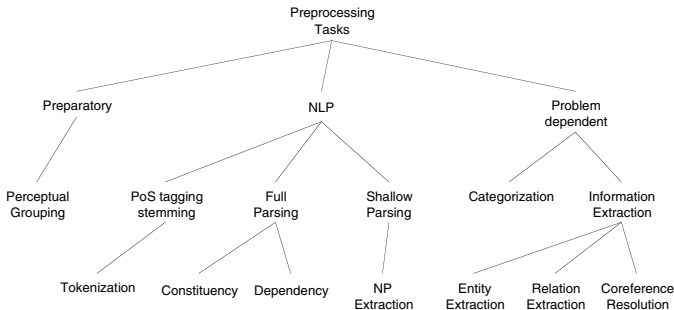
- Pour les datamineurs, la fouille de texte se caractérise par les propriétés suivantes :
 - ① les données (= textes) ont une structure implicite (= la structure linguistique);
 - ② cette structure est flexible car individuelle pour chaque locuteur, mais aussi partagée par tous les locuteurs de la langue;
 - ③ elle est complexe et multicouches : la linguistique définit au moins cinq couches d'étude de la langue :
 - ① phonologie / graphématique,
 - ② morphologie,
 - ③ syntaxe,
 - ④ sémantique,
 - ⑤ pragmatique;
- à l'issue des traitements linguistiques on obtient des données, que l'on peut ensuite fouiller comme d'ordinaire;
- c'est donc la phase de « prétraitement » qui est importante.



Architecture d'un système de fouille de texte



- Le *prétraitement* consiste à traiter le corpus pour obtenir des *représentations de documents* ;
- pour chaque document, une telle représentation peut être simpliste (quelques termes ou concepts dans d'une base de connaissances) ou complexe (le résultat d'une analyse linguistique plus ou moins poussée du texte)



UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- La **tokenisation** consiste à couper le texte en **tokens** (que l'on a va éventuellement classer en tokens de mot, de blanc, de ponctuation, etc.);
- un **token** peut être, selon le besoin, un caractère, une syllabe, un mot, une phrase, un paragraphe, une section, etc.
- pour les caractères : Unicode Standard Annex #29;
- pb avec écritures qui ne séparent pas les mots (chinois, thaï, khmer, mais aussi partiellement arabe, etc.);
- TRÈS IMPORTANT : **MÉTHODE GÉNÉRALE** quand on bloque sur une couche on se sert des couches supérieures;
- on est à la base de l'édifice, si on se trompe maintenant toutes les analyses ultérieures vont se planter !

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- La **phonétique** décrit tous les sons que la bouche de l'humain peut produire ;
- la **phonologie** définit des classes d'équivalence de sons pour une langue donnée, on les appelle **phonèmes** ;
- exemple : en japonais Larry = Rally = ラリ, il s'agit d'un seul phonème la/ra/ラ et d'un seul phonème li/ri/リ ;
- les phonèmes sont validés par la **méthode des paires distinctives** : batte/patte, pif/pouf, pif/pic ;
- la **phonologie dérivationnelle** considère qu'il existe des formes abstraites des mots qui deviennent concrètes à travers des dérivations (→ Chomsky) : in + mature → immature ;
- la **graphématique** étudie les unités élémentaires de l'écriture et leurs relations avec les phonèmes.

UVF3B403 MS

IABDA

Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- La morphologie étudie les *morphèmes* qui sont les briques dont on se sert pour faire des phrases. Leur régularité et leurs oppositions forment une architecture mentale qui nous permet d'exprimer la réalité.
- Un morphème est la plus petite unité porteuse de sens qu'il soit possible d'isoler dans un énoncé.
- Les *morphèmes libres* (ou bases lexicales) peuvent être trouvés seuls et on peut en ajouter des nouveaux à une langue : table, marche, souvent, ...
- Les *morphèmes liés* (ou affixes) ne peuvent être seuls et leur liste est (en principe) fermée : tables, marcher**ont**, etc.
- en japonais, les caractères chinois sont des morphèmes libres et les particules (écrits en hiragana) sont des morphèmes liés.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- La *morphologie flexionnelle* étudie les transformations des mots qui ne produisent pas de mot nouveau (marche → marchera, marchons, marchais, etc., table → tables).
- Principales catégories de la morphologie flexionnelle : pour les noms : le nombre (singulier / pluriel), le genre (masculin / féminin / neutre), le cas (nominatif / génitif / accusatif, etc.). Pour les verbes : la personne (1^{re}, 2^e, 3^e), le mode (impératif / subjonctif / indicatif), la voix (active / passive), le temps (présent, futur, etc.). Chaque langue a sa propre configuration :
 - le français a 2 nombres, 2 genres, aucun cas, le verbe s'accorde avec le nombre du nom ;
 - l'allemand a 2 nombres, 3 genres, 4 cas (n,g,d,a) ;
 - le grec a 3 nombres, 3 genres, 5 cas (n,g,d,a,v) ;
 - le russe a 3 nombres, 3 genres, 6 cas (n,g,d,a,l,i) ;
 - l'arabe a 3 nombres, 2 genres, 3 cas (n,g,a), le verbe dépend aussi du genre du nom ;
 - le japonais n'a ni personne, ni nombre, ni genre, ni cas, tout est indiqué à l'aide de particules ;
 - etc.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- La *morphologie dérivationnelle* produit des nouveaux mots :
table → tablette, marche → démarche, mange → mangeable,
etc.
- Mécanismes dérivationnels : suffixation (manger → mangeable), préfixation (tirer → retirer), composition (canapé-lit, casse-noisettes), siglaison (CNRS, ovni), troncation (sympa, métro, psy), acronymie (franglais, alcootest), reduplication (fefille, pépère).
- En arabe tout mot peut être obtenu à partir d'une *racine sémitique*, le plus souvent trilittère. Par exemple : $\sqrt{\text{ك ت ب}}$
donne كتاب (livre) كتب (livres) كاتب (écrivain) مكتبة (bureau)
يكتب (il écrit), etc.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- Dans un texte on trouve des **occurrences** d'un mot. Quand ces occurrences sont des formes du même mot, on dira qu'il s'agit du même **lexème** (ou unité lexicale). Les unités lexicales utilisées par une personne forment son **vocabulaire**, celles d'une communauté linguistique forment le **lexique** d'une langue.
- On étudie les unités lexicales dans des collections de documents qui se ressemblent par leur origine, que l'on appelle des **corpus**.

POS tagging + stemming : les tags

- La **désuffixation** (stemming) efface les terminaisons d'un mot (par exemple : le radical de « porteront » est port). Cela n'a rien de linguistique mais peut être pratique.
- La **lemmatisation** consiste à trouver la forme canonique d'un mot (par exemple : celle de « porteront » est porter : l'infinitif).
- L'**étiquetage grammatical** (POS tagging) est le processus qui consiste à associer aux mots d'un texte leur fonction grammaticale (sous forme de **tags**).
- Voici la liste de tags du logiciel libre *TreeTagger* <http://www.ims.uni-stuttgart.de/projekte/corplex/TreeTagger/>

| | |
|---------|---|
| ABR | abréviation (CNRS, ENST, ...) |
| ADJ | adjectif (long, court, gros, ...) |
| ADV | adverbe (heureusement, dûment,...) |
| DET:ART | article (le, la, les) |
| DET:POS | pronom possessif (ma, ta, ...) |
| INT | interjection (Ah! Holà! Boum! Atchoum! Berk! ...) |
| KON | conjonction (Mais où e(s)t donc Ornicar? quand) |

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

| | |
|---------|---|
| NAM | nom propre (Noé, Meyer, ...) |
| NOM | nom (<i>cheval, gâteau, espoir, ...</i>) |
| NUM | numéral (15, trois, XXIX, ...) |
| PRO | pronom |
| PRO:DEM | pronom démonstratif (<i>celui, celui-ci, ...</i>) |
| PRO:IND | pronom indéfini (<i>personne, plusieurs, chacun, ...</i>) |
| PRO:PER | pronom personnel (<i>je, me, moi, tu, te, toi, ...</i>) |
| PRO:POS | pronom possessif (<i>mien, tien, ...</i>) |
| PRO:REL | pronom relatif (<i>qui, lequel, quiconque, ...</i>) |
| PRP | préposition (<i>de, à, pour, en, dans, avec, sur, par, ...</i>) |
| PRP:det | préposition plus article (<i>au, du, aux, des</i>) |
| PUN | ponctuation |
| PUN:cit | ponctuation citation (guillemets) |
| SENT | tag de phrase |
| SYM | symbole |

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

| | |
|----------|--|
| VER:cond | verbe au conditionnel (<i>j'irais</i>) |
| VER:futu | verbe au futur (<i>j'irai</i>) |
| VER:impe | verbe à l'imperatif (<i>va</i>) |
| VER:impf | verbe à l'imparfait (<i>j'allais</i>) |
| VER:infi | verbe à l'infinitif (<i>aller</i>) |
| VER:pper | verbe au participe passé (<i>allé</i>) |
| VER:ppre | verbe au participe présent (<i>allant</i>) |
| VER:pres | verbe au présent (<i>je vais</i>) |
| VER:simp | verbe au passé simple (<i>j'allai</i>) |
| VER:subi | verbe à l'imparfait du subjonctif (<i>j'allasse</i>) |
| VER:subp | verbe au présent du subjonctif (<i>j'aille</i>) |

Autres tags fréquemment utilisés (outre-atlantique) : le [Penn Treebank Tagset](http://www.mozart-oz.org/mogul/doc/lager/brill-tagger/penn.html)

Étiquetage : exemple

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

*Longtemps, je me suis couché de bonne heure. Parfois, à peine ma
bougie éteinte, mes yeux se fermaient si vite que je n'avais pas le
temps de me dire : « Je m'endors. »*

donne :

| | | |
|-----------|----------|-------------------------|
| Longtemps | ADV | longtemps |
| , | PUN | , |
| je | PRO:PER | je |
| me | PRO:PER | me |
| suis | VER:pres | sui vre lê tre |
| couché | VER:pper | coucher |
| de | PRP | de |
| bonne | ADJ | bon |
| heure | NOM | heure |
| . | SENT | . |
| Parfois | ADV | parfois |
| , | PUN | , |
| à | PRP | à |
| peine | NOM | peine |
| ma | DET:POS | mon |

Étiquetage : exemple

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

| | | |
|-----------|----------|----------|
| bougie | NOM | bougie |
| éteinte | VER:pper | éteindre |
| , | PUN | , |
| mes | DET:POS | mon |
| yeux | NOM | œil |
| se | PRO:PER | se |
| fermaient | VER:impf | fermer |
| si | ADV | si |
| vite | ADV | vite |
| que | KON | que |
| je | PRO:PER | je |
| n' | ADV | ne |
| avais | VER:impf | avoir |

| | | |
|--------|----------|----------|
| pas | ADV | pas |
| le | DET:ART | le |
| temps | NOM | temps |
| de | PRP | de |
| me | PRO:PER | me |
| dire | VER:infi | dire |
| : | PUN | : |
| « | PUN:cit | « |
| Je | PRO:PER | je |
| m' | PRO:PER | me |
| endors | VER:pres | endormir |
| . | SENT | . |
| » | PUN:cit | » |

Comment fonctionnent les tagueurs ?

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- À partir d'un corpus tagué, on calcule les fréquences des correspondances mot+tag. Si elles sont trop rares (inf. à 1%) on les ignore.
- (Un *n-gramme* est un n -uplet d'objets consécutifs, les objets pouvant être des caractères, des syllabes, des mots, des termes, des phrases, etc.)
- Une première approche : on utilise des n -grammes de mots (souvent $n = 2$ ou 3) et on émet une hypothèse markovienne : le tag d'un mot ne dépend que de celui du précédent, et ce comportement ne varie pas dans le temps.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Soit $X = (X_1, \dots, X_T)$ une suite de v.a. à valeurs dans un espace fini $S = \{s_1, \dots, s_N\}$. Les *propriétés de Markov* sont :
 - 1 l'horizon limité : $P(X_{t+1} = s_k \mid X_1 \dots X_t) = P(X_{t+1} = s_k \mid X_t)$,
 - 2 la stationnarité : $P(X_{t+1} = s_k \mid X_1 \dots X_t) = P(X_2 = s_k \mid X_1)$;
- X est alors une *chaîne de Markov* et on peut la considérer comme un automate dont S sont les états, et dont les transitions de s_i à s_j ont des probabilités $a_{ij} = P(X_{t+1} = s_j \mid X_t = s_i)$ ($a_{*,*}$ est appelée *matrice stochastique*);
- les probabilités des états initiaux s_i sont $\pi_i = P(X_1 = s_i)$ (avec $\sum \pi_i = 1$);
- la probabilité d'avoir une suite avérée d'états X est

$$P(X_1, \dots, X_T) = \pi_{X_1} \prod_{t=1}^{T-1} a_{X_t X_{t+1}}$$

- on appelle cela un *modèle de Markov visible*.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Quand on ne connaît pas la suite d'états mais uniquement une fonction aléatoire de ceux-ci, on parle de *modèle de Markov caché*.
- On procède par apprentissage.
- Soient w_i les mots (resp. t_i les tags) de la phrase à taguer, w^ℓ les mots du lexique (resp. t^j les tags disponibles), $C(w^\ell)$ (resp. $C(t^j)$) le nb. d'occ. de w^ℓ (resp. t^j) dans le corpus d'entraînement, $C(t^j, t^k)$ (resp. $C(w^\ell : t^j)$) le nb. d'occ. de t^j suivi de t^k (resp. du mot w^ℓ avec le tag t^j),
 - alors $\hat{P}(t^k | t^j) = \frac{C(t^j, t^k)}{C(t^j)}$ est le MV (maximum de vraisemblance) que le tag t^k suive le tag t^j ,
 - $\hat{P}(w^\ell | t^j) = \frac{C(w^\ell : t^j)}{C(t^j)}$ est le MV que le mot w^ℓ survienne dans le corpus avec le tag t^j ,
 - alors pour obtenir la suite optimale de tags :

$$\hat{t}_{1,n} = \arg \max_{t_{1,n}} \prod_{i=1}^n \hat{P}(w_i | t_i) \hat{P}(t_i | t_{i-1}).$$

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Les modèles de Markov cachés sont essentiellement des tagueurs à bigrammes. Dans une phrase du type
la jeune fille la porte
impossible de dire à partir de « la » si « porte » est nom ou verbe.
- [Church 1988] décrit un tagueur à trigrammes.
- Dans un tagueur à trigrammes, on peut dire que la probabilité pour une suite de n mot $w_{1,n}$ d'avoir les tags $t_{1,n}$ peut être calculée de manière récursive par

$$p(w_{1,n}, t_{1,n}) := p(t_n | t_{n-2} t_{n-1}) p(w_n | t_n) p(w_{1,n-1}, t_{1,n-1}).$$

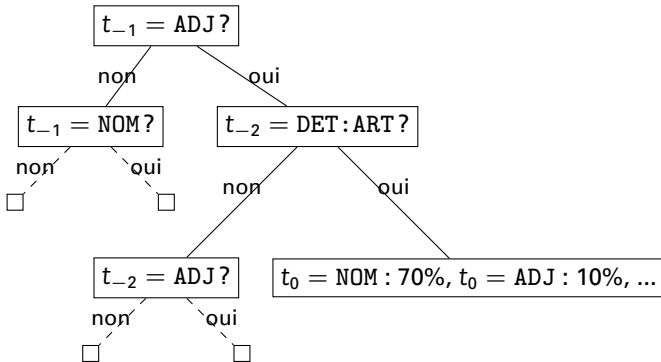
- On peut donc estimer $p(t_n | t_{n-2} t_{n-1})$, en suivant le principe du MV, par

$$\frac{\text{fréquence du trigramme } t_{n-2} t_{n-1} t_n}{\text{fréquence du bigramme } t_{n-2} t_{n-1}}.$$

- Problème : ces probas peuvent être très faibles (voire nulles si le trigramme n'apparaît pas dans le corpus d'apprentissage).

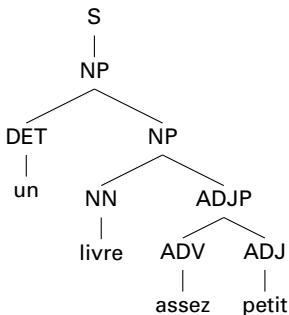
Comment fonctionne *TreeTagger* ?

- Solution adoptée par *TreeTagger* : utiliser un arbre binaire de décision pour déterminer $p(t_n | t_{n-2}t_{n-1})$. Exemple :



- L'apprentissage se fait par ID3 sur des corpus tagués.
- *TreeTagger* dispose d'un lexique de mots avec des probabilités de tags pour le $p(w_n | t_n)$. Si le mot ne s'y trouve pas, il se sert d'expressions régulières sur les suffixes.

- Une *phrase* est une suite de mots qui obéit, dans le système d'une langue donnée, à un ordre et à des dépendances donnés.
- Une phrase se subdivise en *syntagmes* (ou *groupes*) qui sont formés de mots :



- Chaque groupe a une *tête* qui lui donne son nom. On trouve ainsi des groupes nominal NP (le petit chat noir), verbal VP (pense à mes amis), prépositionnel PP (près de la porte), adjectival ADJP (très fier de son travail), ADVP adverbial (tout aussi simplement), etc.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- De manière similaire aux langages formels, l'arbre syntaxique décrit la réalisation d'un mot du langage à partir d'un certain nombre de productions.
- On peut donc décrire la syntaxe d'un langage en décrivant les règles de production :

$$S \rightarrow NP$$
$$NP \rightarrow DET\ NP$$
$$NP \rightarrow NN\ ADJP$$
$$ADJP \rightarrow ADV\ ADJ$$

etc.

- Les feuilles de cet arbre sont les mots. Les parents des feuilles sont les tags de partie du discours (en informatique : les tags POS).

UVF3B403 MS

IABDA

Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- La syntaxe agit sur la morphologie. Ainsi, en français, appliquer la règle $NP \rightarrow DET\ ADJ\ NN\ ADJ$ implique que les adjectifs s'accordent en genre (et en nombre) avec le nom: « des petits livres rouges ».
- Certains traits des verbes interdisent certaines productions, par exemple lorsqu'un verbe est intransitif la production $VP \rightarrow VBZ\ NP$ (la fille dort la pomme) est impossible (sauf en poésie!).
- Autre exemple : « prendre » est toujours suivi d'un COD, sauf dans l'expression « la mayonnaise prend »...
- Les phrases du monde réel sont souvent plus complexes que ces quelques exemples (cf. transparent suivant). Les phrases peuvent être *enchâssées* (Elle sait que Pierre est parti), *coordonnées* (j'ai cherché les clés mais je ne les ai pas trouvées), *transformées* (transformations négative, impérative, interrogative, passive, et leurs combinaisons), etc.

La phrase la plus longue de Proust (243 mots)

UVF3B403 MS

IABDA

Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

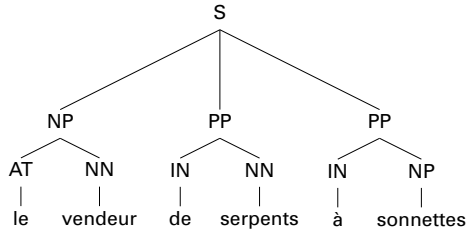
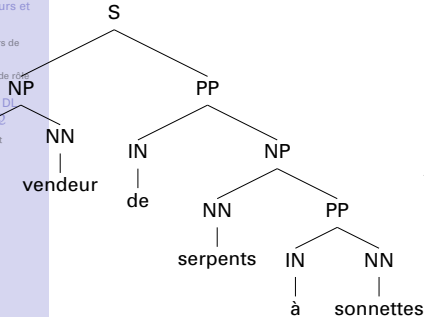
Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

Mais au lieu de la simplicité, c'est le faste que je mettais au plus haut rang, si, après que j'avais forcé Françoise, qui n'en pouvait plus et disait que les jambes « lui rentraient », à faire les cent pas pendant une heure, je voyais enfin, débouchant de l'allée qui vient de la Porte Dauphine — image pour moi d'un prestige royal, d'une arrivée souveraine telle qu'aucune reine véritable n'a pu m'en donner l'impression dans la suite, parce que j'avais de leur pouvoir une notion moins vague et plus expérimentale, — emportée par le vol de deux chevaux ardents, minces et contournés comme on en voit dans les dessins de Constantin Guys, portant établi sur son siège un énorme cocher fourré comme un cosaque, à côté d'un petit groom rappelant le « tigre » de « feu Baudenord », je voyais — ou plutôt je sentais imprimer sa forme dans mon cœur par une nette et épuisante blessure — une incomparable victoria, à dessein un peu haute et laissant passer à travers son luxe « dernier cri » des allusions aux formes anciennes, au fond de laquelle reposait avec abandon M^{me} Swann, ses cheveux maintenant blonds avec une seule mèche grise ceints d'un mince bandeau de fleurs, le plus souvent des violettes, d'où descendaient de longs voiles, à la main une ombrelle mauve, aux lèvres un sourire ambigu où je ne voyais que la bienveillance d'une Majesté et où il y avait surtout la provocation de la cocotte, et qu'elle inclinait avec douceur sur les personnes qui la saluaient.

- Des logiciels appelés *analyseurs syntaxiques* peuvent calculer le (ou les) arbre syntaxique(s) d'une phrase. Parfois il y a plus d'un arbre, correspondant à des sémantiques différentes. Pour désambiguïser on passe aux niveaux supérieurs (sémantique, pragmatique) :



- Autre exemple illustre : la jeune porte le voile.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Pour décrire ces arbres on peut utiliser des *grammaires formelles hors-contexte stochastiques* (SCFG). Ce sont des grammaires formelles hors-contexte (type 2 dans la hiérarchie de Chomsky) avec des probabilités sur les règles de production $P(N^i \rightarrow \zeta^j \mid N^i)$ où N^i est un non-terminal et ζ^j une suite de terminaux et de non-terminaux (avec $\sum_j P(N^i \rightarrow \zeta^j \mid N^i) = 1$ pour tout i).
- La probabilité d'une phrase par rapport à la grammaire G est $P(w_{1m}) = \sum_t P(w_{1m}, t)$ pour tout arbre t dont les feuilles sont w_{1m} .
- On note N_{kl}^j la branche sous N^j dont les feuilles sont les mots w_k, \dots, w_l .
- On pose trois conditions :
 - 1 l'invariance spatiale : $P(N_{k(k+c)}^j \rightarrow \zeta)$ reste la même qq soit k ;
 - 2 hors-contexte : $P(N_{kl}^j \rightarrow \zeta \mid \text{mots en dehors de } w_{kl}) = P(N_{kl}^j \rightarrow \zeta)$;
 - 3 hors-ancêtres : $P(N_{kl}^j \rightarrow \zeta \mid \text{nœuds ancêtres de } N_{kl}^j) = P(N_{kl}^j \rightarrow \zeta)$.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

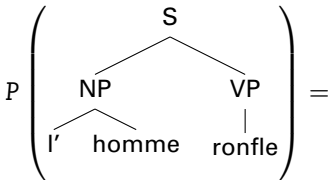
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

Petit exemple d'application des conditions :



$$= P(S \rightarrow NP VP, NP \rightarrow I'homme, VP \rightarrow ronfle)$$

$$= P(S \rightarrow NP VP) \cdot P(NP \rightarrow I'homme \mid S \rightarrow NP VP)$$

$$\cdot P(VP \rightarrow ronfle \mid S \rightarrow NP VP, NP \rightarrow I'homme)$$

$$\stackrel{\text{cond. 2}}{=} P(S \rightarrow NP VP) \cdot P(NP \rightarrow I'homme) \cdot P(VP \rightarrow ronfle)$$

où la cond. 1 nous a permis d'utiliser des expressions générales
alors qu'elles sont extraites d'un arbre spécifique.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Pour faire des calculs on passe à la *forme normale de Chomsky* : on n'admet que des règles du type

① $N^i \rightarrow N^j N^k,$

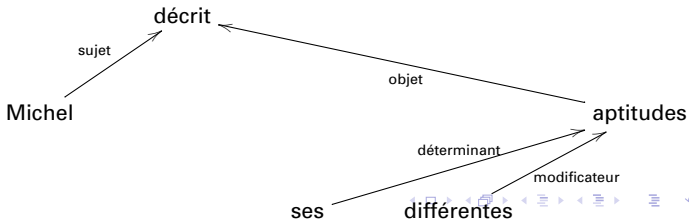
② $N^i \rightarrow w^j.$

donc, pour n non-terminaux et V terminaux, on a $n^3 + nV$ règles possibles.

- On calcule la probabilité d'une phrase selon une grammaire G en utilisant des techniques similaires à celles des HMM.
- Idem pour le calcul de l'arbre le plus probable pour une phrase. L'algorithme de Viterbi nous évite de parcourir tous les arbres possibles (complexité exponentielle).
- Les méthodes d'*inférence grammaticale* (un domaine très actif) nous permettent d'identifier les règles qui produisent un ensemble de phrases prétaguées, et leurs probabilités de manière à coller au plus au corpus d'apprentissage.

- Les SCFG ne tiennent pas compte des mots et de leurs propriétés. Une phrase du type « L'homme dort la lune » sera considérée comme correcte puisque la SCFG ignore l'intransitivité du verbe dormir.
- (Parenthèse : Ainsi, le Lefff [Sagot 2010] nous donne les infos suivantes sur la forme « dort » :
dort [pred='dormir____1<Suj:cln|sn>',@pers,cat=v,@P3s]
alors que pour le verbe « prendre » il nous aurait dit :
prend
[pred='prendre____2<Suj:cln|scompl|sinf|sn,Obj:(cla|sn),
Objà:(cld|à-sn)>',@pers,cat=v,@P3s])
- D'autre part, il paraît que l'hypothèse d'invariance spatiale est fausse dans la réalité : par exemple, les noms propres apparaissent bcp plus souvent en tant que sujets qu'en tant qu'objets.

- La syntaxe que nous avons apprise à l'école se base sur une grammaire de constituants : le sujet (ou le GN), le verbe (ou le GV), l'objet, etc.
- Il existe une autre manière de décrire la structure d'une phrase : la *grammaire des dépendances* de Tesnière.
- On considère une phrase comme un mot (la *tête* : le plus souvent, le verbe) auquel sont attachés des *modificateurs*, qui peuvent avoir des modificateurs à leur tour. Le concept fondamental est celui de *relation* entre les mots, relation fléchée entre *dépendant* et *gouverneur*.
- Exemple d'arbre de dépendances :



Autre exemple d'arbre de dépendances

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous

(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et restrictions

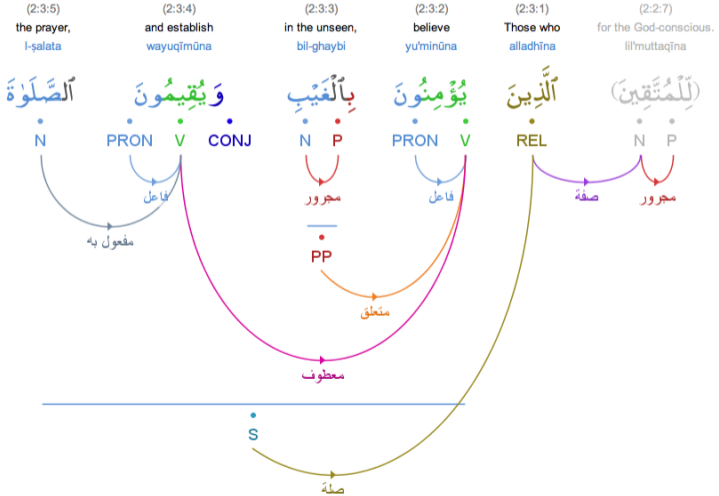
Constructeurs de concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et exemples

Requêtes

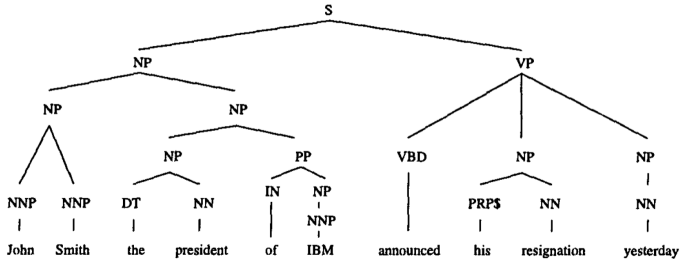


15/06/10 13:44

Source :

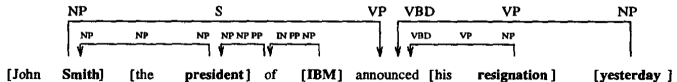
http://corpus.guran.com/treebank_isp?chapter=2&verse=3&token=1

- [Collins 1996, 1997] présente un modèle statistique et un parseur lexicalisé, basés sur la grammaire des dépendances.
- Le modèle statistique calcule, pour tout arbre T de la phrase S , la probabilité $P(T \mid S)$. Pour cela, S est réduit en $\bar{S} = (B, D)$ (bases et dépendances).
- Exemple :



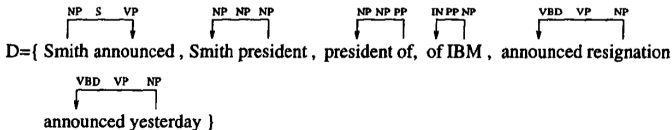
27/06/10 11:20

- Voici les dépendances qui en découlent :



- et voici les ensembles B et D :

$B = \{ [John\ Smith], [the\ president], [IBM], [his\ resignation], [yesterday] \}$



27/06/10 11:21

- Pour réduire l'arbre, il va :
 - pour chaque nœud intermédiaire trouver le **head-child** par un arbre de décisions [Jelinek et al. 1994];
 - faire remonter les **head-child** aussi haut que possible;
 - extraire les dépendances.
- Chaque constituant à n enfants produit $n - 1$ dépendances.
- Ayant mis l'arbre à plat, le parseur fait de l'apprentissage sur des triplets des dépendances.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- Le *shallow parsing* (ou *chunking*) est la version allégée de l'analyse syntaxique : au lieu de trouver un arbre syntaxique complet, on se contente d'une structure superficielle. Il convient parfaitement à l'extraction d'information ou à la classification de documents.
- [Abney 1994] définit la notion de *chunk* et la justifie en se basant sur des études psychologiques.
- Pour lui, un chunk est tout mot non-grammatical w accompagné de mots grammaticaux qui l'entourent, *sauf* quand w se trouve entre un autre mot w' et un mot grammatical dépendant de w' . Ainsi, alors que « big », « proud » et « man » sont des noms, « a big proud man » est un seul chunk, puisque « a » dépend de « man ».

- [Sha & Pereira 2003] définissent les chunks comme étant les noyaux non-récursifs de différents groupes de la phrase.
- Les chunks étant non-récursifs, le découpage de la phrase est toujours plat. Exemple :

[Indexing]_N [for the most part]_N [has involved simply buying]_V
[and then holding]_V [stocks]_N [in the correct mix]_N [to mirror]_V
[a stock market barometer]_N.

- L'*Illinois Chunker*

http://cogcomp.cs.illinois.edu/page/software_view/13
trouve des chunks de 10 types (ADJP, ADVP, CONJP, INTJ, LST,
NP, PP, PRT, SBAR, VP):

[NP Jack and Jill] [VP went] [ADVP up] [NP the hill] [VP to
fetch] [NP a pail] [PP of] [NP water]

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Le problème de base (on se restreint aux chunks nominaux) consiste à associer à chaque mot un « label »
 - 1 O s'il est à l'extérieur de tout chunk ;
 - 2 B s'il est le premier mot d'un chunk ;
 - 3 C s'il appartient à un chunk, sans en être le premier mot.
- Pour faire cela, [Ramshaw & Marcus 1995] utilisent des règles de transformations faisant intervenir jusqu'à 3 mots ou tags POS à gauche et à droite d'un mot donné.
- Entre 1992 et 2001, plusieurs techniques ont été utilisées : les HMM, les MEMM, Winnow, AdaBoost, les SVM, le perceptron généralisé.
- [Sha & Pereira 2003] utilisent des CRF pour attaquer ce problème. Pour un mot w_i ils font intervenir w_{i-2}, \dots, w_{i+2} , t_{i-2}, \dots, t_{i+2} (les tags POS), y_i et y_{i+1} (les labels de chunk).

- La sémantique est l'étude (scientifique) de la signification.
- Dans un processus d'abstraction à partir des choses qui nous entourent on arrive aux *concepts* et aux *relations* entre les concepts.
- Au-delà des considérations philosophique, on peut se demander comment décrire, classer, calculer les concepts ?
- Une méthode classique consiste à énumérer leur *traits* (= catégorie + valeur) :

| français | allemand | échange | donner qqch | recevoir qqch | pour argent | permanent |
|-----------|-----------|---------|-------------|---------------|-------------|-----------|
| acheter | kaufen | + | - | + | + | + |
| vendre | verkaufen | + | + | - | + | + |
| prêter | ausleihen | + | + | - | - | - |
| emprunter | leihen | + | - | + | - | - |
| louer | vermieten | + | + | - | + | - |
| louer | mieten | + | - | + | + | - |

Cette méthode a inspiré la discipline FCA (analyse de concepts formels).

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

On a un certain nombre de phénomènes :

- la **polysémie** : le même mot a plusieurs sens, liés entre eux (exemple : la construction du pont, cet édifice est une belle construction) ;
- l'**homonymie** : plusieurs sens sans rapport entre eux (exemple : lire un livre, une livre sterling) ;
- l'**hyponymie** : quand on a une relation d'implication dans un sens (exemple : le concept de pomme est un hyponyme du concept de fruit, puisque toute pomme est fruit, mais tout fruit n'est pas pomme) ;
- l'**hypéronymie** : l'inverse de l'**hyponymie** ;
- la **synonymie** : quand les implications vont dans les deux sens (exemple : soulier et chaussure). Attention : il peut y avoir des différences stylistiques ou historiques, on n'a jamais de synonymie parfaite (on parle aussi de **cohyponymie**) ;
- la **méronymie** : quand le premier est une partie du second (exemple : le doigt est une partie de la main qui est une partie du corps).

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Voici comment cette théorie définit les notions de *contexte*, *concept*, *extension* et *intension*.
- Un *contexte* est un triplet (G, M, I) où G (objets) et M (traits ou attributs) sont des ensembles et $I \subseteq G \times M$. $(g, m) \in I$ signifie que « g possède l'attribut m »;
- pour $G_1 \subset G$, $M_1 \subset M$ les dérivés G_1^I et M_1^I sont définis par

$$G_1^I = \{m \in M \mid \forall g \in G_1, (g, m) \in I\}$$

$$M_1^I = \{g \in G \mid \forall m \in M_1, (g, m) \in I\}$$

autrement dit : G_1^I est l'ensemble des attributs partagés par tous les éléments de G_1 ; M_1^I est l'ensemble des objets possédant tous les attributs de M_1 ;

- un *concept* (X, N) de (G, M, I) est un couple $X \subseteq G, N \subseteq M$ tel que $X^I = N$ et $N^I = X$. On dit que X est l'*extension* du concept et N son *intension*.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- À partir d'un objet g on peut obtenir son **concept d'objet** $(\{g\}^{\text{II}}, \{g\}^{\text{I}})$;
- à partir d'un attribut m on peut obtenir son **concept d'attribut** $(\{m\}^{\text{I}}, \{m\}^{\text{II}})$;
- on a $X \subseteq X^{\text{II}}, N \subseteq N^{\text{II}}, X^{\text{I}} = X^{\text{III}}, N^{\text{I}} = N^{\text{III}}$;
- on a un ordre partiel $(X_1, N_1) \leq (X_2, N_2)$ induit par l'inclusion dans G (\Leftrightarrow l'inverse de l'inclusion dans N) : plus un concept est grand (= large), plus il contient des objets, moins il contient des attributs ;
- le concept d'objet de g est le plus petit concept dont l'extension contient g / le concept d'attribut de m est le plus grand concept dont l'intension contient m ;
- muni de l'ordre \leq , un contexte a une structure de **treillis** (plus précisément, de **treillis de Galois** pour la relation binaire I).

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- En informatique, pour associer des mots à des concepts et définir des relations entre eux, on utilise les *ontologies*.
- Une ontologie comporte des *concepts*, des *relations* entre les concepts, des *attributs*, une *hiérarchie de concepts*, une *hiérarchie de relations* définie à partir de la hiérarchie de concepts, et un ensemble de types (chaîne, entier, etc.).

Mathématiquement parlant :

- 1 Une ontologie $\mathcal{O} := (C, \leq_C, R, \sigma_R, \leq_R, A, \sigma_A, T)$
 - 2 où C, R, A, T sont les concepts, relations, attributs, types ;
 - 3 $\sigma_R : R \rightarrow C^+, \sigma_A : A \rightarrow C \times T$ sont les signatures des relations et des attributs ;
 - 4 \leq_C munit C d'une structure de treillis ;
 - 5 \leq_R est un ordre partiel de R , défini de la manière suivante :
 $r_1 \leq_R r_2 \Rightarrow |\sigma_R(r_1)| = |\sigma_R(r_2)|$ et $\pi_i(\sigma_R(r_1)) \leq_C \pi_i(\sigma_R(r_2))$ pour tout i (ici π_i est la i -ème projection).
- Exemple : concepts rivière, ville, relation `traverse()`, signature $\sigma_R(\text{traverse}) = (\text{rivière}, \text{ville})$.

UVF3B403 MS

IABDA

Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- À l'aide du λ -calcul on définit des schémas d'axiomes, et on en munit l'ontologie. Exemple : l'axiome $\lambda P, Q. \text{disjoint}(P, Q)$, peut être décrit dans le formalisme de logique du premier ordre $\lambda P, Q. \forall x (P(x) \rightarrow \neg Q(x))$ et donc, par exemple, si x est rivière il ne peut pas être montagne.
- Si on veut peupler une ontologie à partir du texte, il faut associer des mots aux concepts. On définit un **lexique** \mathcal{L} pour une ontologie \mathcal{O} comme étant $\mathcal{L} := (S_C, S_R, S_A, \text{Ref}_C, \text{Ref}_R, \text{Ref}_A)$ où S_C, S_R, S_A sont des lexèmes pour les concepts, relations et attributs, et $\text{Ref}_* \subseteq S_* \times *$ des **références lexicales** pour les concepts, relations et attributs.
- Ainsi, par exemple, quand on a des synonymes (voiture, auto, bagnole, caisse) on pourra dire que $\text{Ref}_C^{-1}(\text{voiture}) = \{\text{voiture, auto, bagnole, caisse}\}$, où voiture dénote le concept de voiture.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Une *base de connaissances* \mathcal{K} pour une ontologie \mathcal{O} est un quadruplet $(I, \iota_C, \iota_R, \iota_A)$ où I est un ensemble d'*instances* (ou *objets*) et ι_* des fonctions des instantiation :

$$\iota_C : C \rightarrow 2^I, \iota_R : R \rightarrow 2^{I^+}, \iota_A : A \rightarrow 2^I \times \text{valeurs}(T);$$
- Même question que tout à l'heure : trouver des entités lexicales pour nos instances. Même réponse :
- Un *lexique d'instances* pour une base de connaissances \mathcal{K} est une paire (S_I, R_I) où S_I sont des signes pour les instances et $R_I \subseteq S_I \times I$ est une relation qui associe des références lexicales aux instances.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- Il est temps de profiter de la structure de treillis de l'ontologie pour obtenir plus d'instances (exemple : on sait qu'un labrador est un toutou, et on a Bill, une instance de labrador, il faudrait que ce soit également une instance de toutou), idem pour les relations et les attributs. Si $c \in C$ est un concept, on définit son extension $\llbracket c \rrbracket_{\mathcal{K}}$ dans \mathcal{K} comme $\llbracket c \rrbracket_{\mathcal{K}} \subseteq I$ par la construction récursive suivante :
 - condition initiale : $\llbracket c \rrbracket_{\mathcal{K}} \leftarrow \iota_C(c)$;
 - pour tout $c' <_C c$, $\llbracket c \rrbracket_{\mathcal{K}} \leftarrow \llbracket c \rrbracket_{\mathcal{K}} \cup \llbracket c' \rrbracket_{\mathcal{K}}$.
- On peut également définir un concept ou une relation intensionnellement. Exemple : le concept `nombre_pair` se définit par $\{n \in \mathbb{Z} \mid \exists p \text{ tel que } n = 2p\}$.

- Extraire de la connaissance à partir de textes dans le but d'alimenter une ontologie revient à extraire les couches successives suivantes d'information :
 - 1 des termes (fleuve, rivière, pays, ville, capitale, ...) et les entités nommées (Brest, IMT Atlantique, Laury Thilleman, etc.);
 - 2 des synonymes { fleuve, rivière, cours d'eau, ... };
 - 3 des concepts *rivière*, *capitale*, *ville*, ... (un concept est une paire de définitions intentionnelle et extensionnelle, ainsi qu'un lexème référent);
 - 4 des hiérarchies de concepts $\text{capitale} \leq_C \text{ville}$;
 - 5 des relations entre concepts $\text{traverse}(\text{rivière}, \text{pays})$;
 - 6 des hiérarchies de relations $\text{est_capitale} \leq_R \text{est_situé_dans}$;
 - 7 des schémas d'axiomes $\text{disjoint}(\text{rivière}, \text{ville})$;
 - 8 des axiomes logiques généraux : $\forall x(\text{pays}(x) \rightarrow \exists y \text{ capitale_de}(y, x)) \wedge \forall z(\text{capitale_de}(z, x) \rightarrow y = z)$.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Les *logiques de description* sont des variantes de la logique du 1^{er} ordre qui forment des compromis sur deux tableaux :
 - 1 elles ajoutent des nouvelles notations pour améliorer l'expressivité du langage (par ex. la possibilité de dire qu'il existe exactement n éléments ayant telle propriété);
 - 2 elles sont moins puissantes que la logique du 1^{er} ordre, ce qui les rend décidables.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

• On a trois types d'objets :

- 1 des *individus* (ce qui correspond aux constantes de la logique du 1^{er} ordre),
- 2 des *concepts* ou *classes* (des prédicats unaires, dont l'interprétation ensembliste correspond à des ensembles d'individus),
- 3 des *rôle* ou *propriétés* (des prédicats binaires sur les individus) ;

UVF3B403 MS
IABDA

Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- On définit trois types de formules :

- 1 celles de l'**ABox**, « A » comme « assertion », qui décrivent des concepts et des rôles (prédicats unaires et binaires sur des constantes) ;
- 2 celles de la **TBox**, « T » comme « terminologie », qui décrivent des relations entre concepts ;
- 3 celles de la **RBox**, « R » comme « relation », qui décrivent une hiérarchie des rôles, des compositions de rôles et des relations entre rôles, comme l'exclusion mutuelle, la réflexivité, la symétrie, la transitivité, etc.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- Une *assertion de concept* est un prédicat unaire du type `ÉlèveTélécom(xavier)`.
- Une *assertion de rôle* est un prédicat binaire `Père(andré,mathilde)`.
- Dans les logiques de description on ne fait pas l'*hypothèse de l'unicité des noms* : deux individus de même nom peuvent avoir le même référent, on écrira `cloclo ≈ claudFrançois`, et `samsonFrançois ≠ claudFrançois`.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Puisqu'on peut interpréter les *concepts* (prédicats unaires) comme des ensembles, on peut aussi utiliser des relations de théorie d'ensembles :
- Mère notations] \sqsubseteq (inclusion de concept) \sqsubseteq Parent, qui équivaut à $\forall X \text{ Mère}(X) \rightarrow \text{Parent}(X)$;
- Personne notations] \equiv (équivalence de concept) \equiv Humain, qui équivaut à $\forall X \text{ Personne}(X) \leftrightarrow \text{Humain}(X)$.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- La *hiérarchie des concepts* induit une *hiérarchie des rôles* :
- $\text{ParentDe} \sqsubseteq \text{AncêtreDe}$, qui équivaut à
 $\forall X, Y \text{ ParentDe}(X, Y) \rightarrow \text{AncêtreDe}(X, Y)$.
- Une relation binaire peut se combiner avec une autre : sachant que le frère d'un père est un oncle, on peut écrire
- $\text{FrèreDe} \circ \text{ParentDe} \sqsubseteq \text{OncleDe}$, qui équivaut à
 $\forall X, Y, Z \text{ FrèreDe}(X, Y) \wedge \text{ParentDe}(Y, Z) \rightarrow \text{OncleDe}(X, Z)$.
- Enfin, on peut donner des caractéristiques générales des concepts : *Disjoint*(ParentDe, EnfantDe), ce qui équivaut à
 $\forall X, Y \text{ ParentDe}(X, Y) \rightarrow \neg \text{Parent}(Y, X)$.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- Toujours en mimant la théorie des ensembles, on peut écrire notations] \sqcap (intersection de concept) $\text{Mère} \equiv \text{Femme} \sqcap \text{Parent}$, qui équivaut à $\forall X \text{Mère}(X) \leftrightarrow \text{Femme}(X) \wedge \text{Parent}(X)$,
- ou notations] \sqcup (union de concept) $\text{Parent} \equiv \text{Père} \sqcup \text{Mère}$, qui équivaut à $\forall X \text{Parent}(X) \leftrightarrow \text{Mère}(X) \vee \text{Père}(X)$,
- la négation correspond au complémentaire d'un ensemble : $\neg \text{Célibataire}$ est équivalent à Marié .
- L'ensemble complet correspond à un prédicat qui est toujours vrai, on l'écrit notations] \top (ensemble complet) \top , pour exprimer une partition on écrira $\top \sqsubset \text{Homme} \sqcup \text{Femme}$.
- L'ensemble vide correspond à un prédicat qui est toujours faux, on l'écrit notations] \perp (ensemble vide) \perp , pour exprimer une exclusion mutuelle on écrira $\text{Homme} \sqcap \text{Femme} \sqsubset \perp$.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Imaginons qu'on a un rôle $\text{Parent}(X, Y)$ et que l'on cherche à caractériser les X qui sont parents. Il s'agit donc de dire « je cherche les X pour lesquels $\exists Y$ tel que $\text{Parent}(X, Y)$ », on écrira $\exists \text{Parent}.\top$;
- le \top signifie qu'on prend les X pour lesquels il existe un Y , sans les filtrer davantage.
- Si je cherchais ceux qui sont des parents d'au moins une fille, j'écrirais $\exists \text{Parent}.\text{Fille}$.
- De même, ceux qui n'ont que des filles : $\forall \text{Parent}.\text{Fille}$.
- Mais que signifie alors $\forall \text{Parent}.\top$?

- Pour le savoir, prenons les définitions formelles des sémantiques de ces notations.
- Soit R un rôle, et R^I son interprétation. Rappelons que dans une interprétation ensembliste, R^I devient un ensemble de paires d'éléments $R^I = \{(x^I, y^I), \dots\}$ du domaine Δ . On dira que y^I est un R -*successeur* de x^I si la paire (x^I, y^I) appartient à R^I .
- Soit C un concept (et donc C^I est un sous-ensemble de Δ). Alors l'interprétation de $\exists R.C$ (quelques successeurs) est $\{x^I \mid \text{quelques successeurs de } x^I \text{ sont dans } C^I\}$.
- Et celle de $\forall R.C$ (tous les successeurs) est $\{x^I \mid \text{tous les successeurs de } x^I \text{ sont dans } C^I\}$.
- Donc, quelle est l'interprétation de $\forall \text{Parent}.T$?

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- De même que l'on peut demander au moins un successeur ou tous les successeurs, on peut aussi spécifier « au moins n successeurs » : $\geq nR.C$ (au moins n successeurs), ainsi qu'« au plus n successeurs » : $\leq nR.C$ (au plus n successeurs).
- Une autre notation permet de décrire l'ensemble des éléments pour lesquels un rôle R est réflexif : $\exists R.Self$ (rôle réflexif).
- Une manière de décrire un concept est en donnant explicitement ses membres : au lieu de $Parent(jacques, julie)$, on peut aussi écrire $\{jacques\} \sqsubset \exists Parent.\{julie\}$.
- On note R^- le *rôle inverse* de R , par exemple $Parent^-$ est équivalent à $Enfant$ puisque $\forall X, Y \text{ Parent}(X, Y) \leftrightarrow \text{Enfant}(Y, X)$.
- Enfin, on note U le *rôle universel*, c'est-à-dire celui qui associe chaque élément à chaque autre élément (y compris lui-même).
- Enfin, on note R^+ la *clôture transitive* de R , c'est-à-dire le plus petit rôle transitif contenant R .

- Une *ontologie* est la formalisation d'un domaine de connaissances dans un langage de représentation de connaissances.
- Une *ontologie DL* est une ontologie représentée dans le langage d'une logique de description.
- *SROIQ* est une logique de description spécifique que nous allons décrire dans la suite.
- Si N_C est un concept quelconque, N_R un rôle quelconque et N_I un individu quelconque, alors on définit de manière itérative dans *SROIQ* :
 - ① une *expression de concept* C par $C ::= N_C \mid (C \sqcap C) \mid C \sqcup C \mid \neg C \mid \top \mid \perp \mid \exists R.C \mid \forall R.C \mid \geq n R.C \mid \leq n R.C \mid \exists R.Self \mid \{N_I\}$,
 - ② une *expression de rôle* R par $R ::= U \mid N_R \mid N_R^-$.
- Les axiomes d'une ontologie *SROIQ* sont des types suivants :
 - ① ABox : $C(N_I), R(N_I, N_I), N_I \approx N_I, N_I \not\approx N_I$,
 - ② TBox : $C \sqsubseteq C, C \equiv C$,
 - ③ RBox : $R \sqsubseteq R, R \equiv R, R \circ R \sqsubseteq R, Disjoint(R, R)$.

- Voici les propriétés de la logique de description *SROIQ* :
 - ❶ \mathcal{S} (également appelé \mathcal{ALCR}^+) : (\mathcal{AL}) présence de noms de concepts et de rôles, de \top , du constructeur \sqcap (conjonction), du quantificateur universel, (\mathcal{C}) de la négation de concept, (\mathcal{R}^+) de la clôture transitive ;
 - ❷ \mathcal{R} : inclusion de rôles, réflexivité, irreflexivité, exclusion de rôles ;
 - ❸ \mathcal{O} : possibilité de décrire un concept extensionnellement $(\{N_{I,1}, \dots, N_{I,n}\})$;
 - ❹ \mathcal{I} : possibilité d'avoir des rôles inverses (\mathcal{R}^-) ;
 - ❺ \mathcal{Q} : possibilité d'avoir des restrictions pleinement quantifiées $(\leq n, \geq n)$.
- D'autres logiques de description sont éventuellement dotées des propriétés suivantes :
 - ❻ \mathcal{H} : hiérarchie des rôles $R_1 \sqsubseteq R_2$, \mathcal{H} est plus faible que \mathcal{R} ;
 - ❼ \mathcal{N} : restrictions de cardinalité plus faibles que \mathcal{Q} .

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SR_OIQ

Définitions et
exemples

Requêtes

- La norme **OWL** 2 du consortium **WWW** correspond à la logique *SR_OIQ*, alors que la norme plus faible OWL-DL correspond à *SHOIN*.
- Le logiciel **Protégé** d'édition d'ontologies est basé sur cette dernière.

Donnée \neq Information \neq Connaissance

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Nous utilisons tous les termes « donnée », « information », « connaissance » dans la vie courante.
- Mais ils ont aussi un sens plus technique, utilisé en intelligence artificielle :
- les *données* sont des ensembles ordonnés de nombres, provenant de capteurs ou de programmes de simulation, de générateurs de données aléatoires, etc. ;
- les *informations* sont des données auxquelles on a attaché un sens : on sait ce que représente tel nombre ou tel ensemble d'octets ou telle chaîne de caractères ;
- les *connaissances* sont des motifs ou des tendances qu'on extrait des informations, dans le but de faire des prédictions.
- Exemple : le départ d'un train de la gare de Brest.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Un *domaine de connaissances* est un ensemble de données, informations et connaissances autour d'un thème donné (par exemple : la géographie est un domaine de connaissances).
- Une *base de connaissances* est une structure informatique/mathématique qui permet de stocker des données, informations et connaissances. Point commun avec les SGBD : on peut interagir avec une base de connaissances par le biais de requêtes/réponses.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Dans ce chapitre, qui est un avant-goût du module de logique INF 424, on va représenter une base de connaissances sous forme de graphe.
- On utilisera pour cela un genre de graphe bien particulier : les *graphes conceptuels*.
- Pour faire le lien entre le domaine que l'on veut décrire et la structure de graphe, on définit la notion de *vocabulaire* :

Definition

Un *vocabulaire de graphe conceptuel* est un triplet (T_C, T_R, \mathcal{I}) , où T_C est l'ensemble de *types de concept*, il est partiellement ordonné et a un plus grand élément \top ; T_R est l'ensemble de *symboles de relation*, il est également partiellement ordonné et chacun de ses éléments a une *arité* ≥ 1 ; \mathcal{I} est l'ensemble des *marqueurs d'individu*. On ajoute à \mathcal{I} l'élément $*$ (*marqueur générique*) avec la propriété $\forall i \in \mathcal{I}, i \geq *$. Ces ensembles sont mutuellement disjoints.

Definition

Un *graphe conceptuel de base* sur un vocabulaire \mathcal{V} donné, est la donnée d'un graphe G biparti, de partitions C et R , non-orienté avec éventuellement des arêtes multiples et d'une fonction ℓ , définis de la manière suivante :

- les sommets sont des *concepts* C et des *relations* R
- l'image de $c \in C$ par ℓ est une paire (t, i) avec $t \in T_C$ et $i \in \mathcal{I}$;
- pour $r \in R$, $\ell(r) \in T_R$;
- le degré de $r \in R$ est égal à l'arité de $\ell(r)$:

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

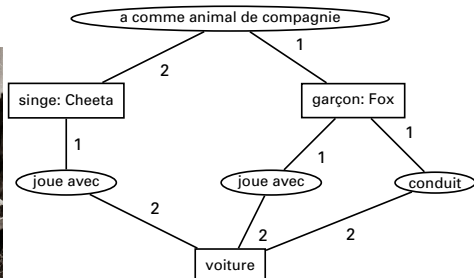
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

Prenons un exemple :



La photo est interprétée par « Fox est un garçon, Cheeta est un singe. Cheeta est l'animal de compagnie de Fox. Ils jouent ensemble avec une voiture-jouet, conduite par Fox. »

Ici, « singe », « garçon » et « voiture » sont des types de concepts, « joue avec » et « conduit » sont des relations d'arité 2 (sujet, COD), « Cheeta » et « Fox » sont des marqueurs d'individu, alors que le sommet du bas est en fait « voiture : * » (on ne note pas le marqueur générique).

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SR OIQ

Définitions et
exemples

Requêtes

Definition

On définit un *ordre des concepts* à partir de l'ordre des types de concept : si $(t, i), (t', i') \in C$, $(t, i) \leq (t', i')$ si et seulement si $t \leq t'$ et $i \leq i'$. Si $t \leq t'$ on dira que t' est un *hypéronyme* de t .

(À noter que $i \leq i'$ ne peut arriver que si $i = i'$ ou si $i' = *$, puisque les i, i' différents de $*$ ne peuvent être comparés.)

Definition

Si G et G' sont des graphes conceptuels sur le même vocabulaire, un *homomorphisme de graphes conceptuels* est un homomorphisme de graphes φ qui envoie C dans C' , R dans R' , $\mathcal{I} \cup \{*\}$ dans $\mathcal{I}' \cup \{*\}$ et qui est tel que $\varphi(c) \leq c$, $\varphi(r) \leq r$ et φ respecte les numéros des arêtes.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

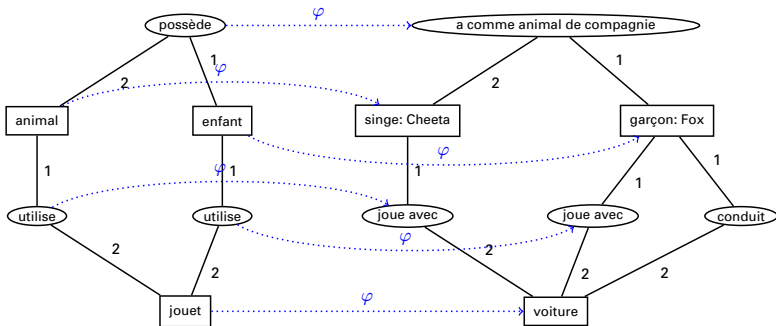
Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes



Definition

Soient G et G' deux graphes conceptuels sur le même vocabulaire. On définit la *relation de subsumption* \succeq de la manière suivante : $G \succeq G'$ ssi il existe un homomorphisme de graphe conceptuel $\varphi: G \rightarrow G'$.

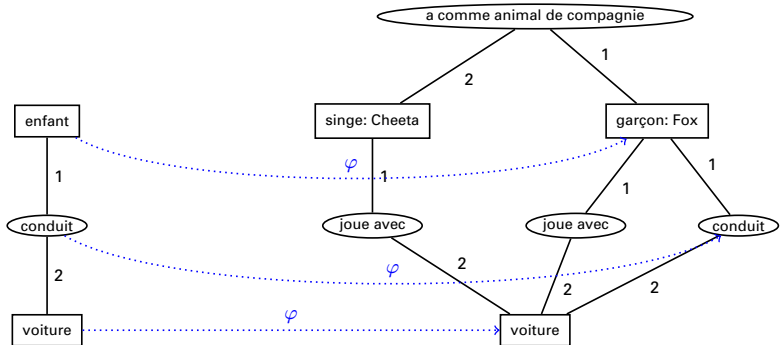
Definition

Soit Q et G des graphes conceptuels sur le même vocabulaire. On dira que Q est une requête acceptée par G si $Q \succeq G$. Les résultats de la requête sont les images $\varphi(Q)$ pour tout homomorphisme $\varphi: Q \rightarrow G$.

On peut imaginer diverses situations où la notion de requête peut être utile :

- (1) on considère que G est un graphe conceptuel qui décrit la réalité et on voudrait savoir si Q est « conforme » à cette réalité (on dira en INF 424 que Q est une « conséquence » de G);
- (2) G décrit l'image ci-dessus, Q correspond à une requête d'image

Premier exemple : « Un enfant conduit une voiture ».



UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

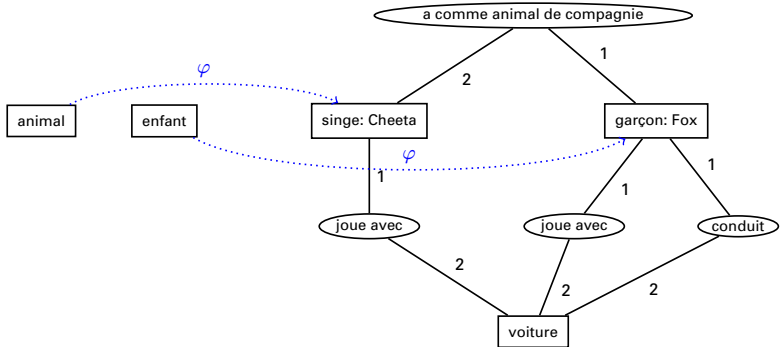
Restrictions de rôle

Ontologies DL
SROIQ

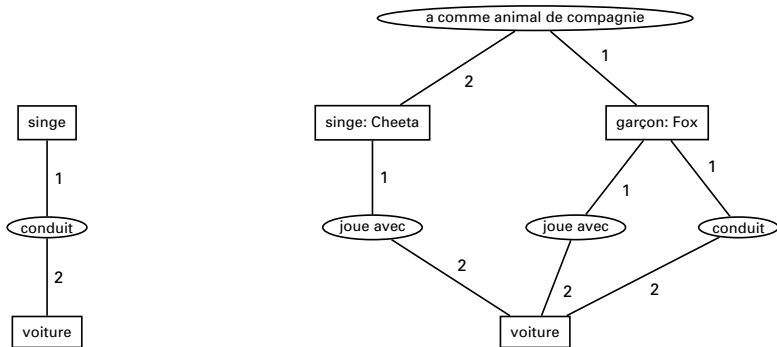
Définitions et
exemples

Requêtes

Deuxième exemple : « Un enfant possède un animal ».

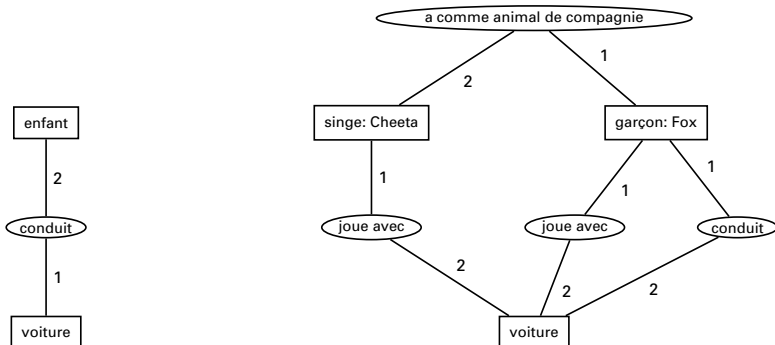


Troisième exemple : « Un singe conduit une voiture ».



Impossible puisque l'arête entre « singe » et « conduit » ne peut être envoyée à « singe : Cheeta »—« conduit » qui n'existe pas. De même, on ne peut pas envoyer « conduit » vers « joue avec » puisqu'elles sont incomparables.

Quatrième exemple : « Une voiture conduit un enfant ».



Impossible puisque l'application qui envoie « enfant » à « garçon : Fox », « conduit » à « conduit » et « voiture » à « voiture », n'est pas un homomorphisme de graphes contextuels puisqu'elle ne respecte pas les numéros des arêtes.

Exemple de problème : la désambiguïsation de mot

- (Officiellement, depuis 1990 ce mot s'écrit « désambiguïsation ».)
- Il s'agit de choisir entre les différents sens d'une instance de mot polysémique (*hier mon fils a volé pour la première fois : voleur ou aviateur?*).
- Première approche [Gale et al. 1992] : **classification bayésienne**. On considère w un mot ambigu, s_* ses sens, c_* les contextes de w dans un corpus, v_* les mots environnants. Pour un c donné, on cherche à maximiser $P(s_i | c)$.

Approche bayésienne : on cherche

$$\arg \max P(s_* | c) = \arg \max [\log P(c | s_*) + \log P(s_*)].$$

Hypothèse « naïve Bayes » : $P(c | s_k) = \prod_{v_j \in c} P(v_j | s_k)$.

On procède par estimation de la vraisemblance maximale :

$$\hat{P}(v_j | s_k) = \frac{\#\{v_j \in \text{contexte}(s_k)\}}{\#\{w\}}, \quad \hat{P}(s_k) = \frac{\#\{s_k\}}{\#\{w\}}$$

dans le corpus d'apprentissage.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Autre approche [Lesk 1986] : se servir des différentes définitions du mot dans un dictionnaire. On prend les mots qui apparaissent dans chacune des définitions d'un mot, et on les compare avec ceux des définitions des mots voisins. [Walker 1987] utilise un thésaurus. [Resnik 1995] utilise les synsets de WordNet.
- WordNet est une base de données lexicale libre. Pour chaque mot, WordNet fournit des synsets (= ensembles de synonymes) et une courte description. Exemple (pour *car*) :
 - car, auto, automobile, machine, motorcar – (4-wheeled motor vehicle; usually propelled by an internal combustion engine; he needs a car to get to work),
 - car, railcar, railway car, railroad car – (a wheeled vehicle adapted to the rails of railroad; three cars had jumped the rails),
 - car, gondola – (car suspended from an airship and carrying personnel and cargo and power plant),
 - car, elevator car – (where passengers ride up and down; the car was on the top floor),
 - cable car, car – (a conveyance for passengers or freight on a cable railway; they took a cable car to the top of the mountain).

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- IE (*Information Extraction*) : recherche d'une information dans un document ou dans un corpus, à ne pas confondre avec IR (*Information Retrieval*) : recherche de documents contenant une certaine information.
- Il s'agit de récupérer des entités, des relations entre entités, des qualificatifs d'entités, ainsi que des structures plus élaborées comme des tableaux et des listes.
- Une *entité* est typiquement une phrase nominale comportant entre un et un petit nombre de mots. Une *entité nommée* est un nom propre, un acronyme, une abréviation, un terme temporel, une expression monétaire ou numérique. En médecine une EN sera un nom de maladie, de protéine, de gène, etc.
- Une *relation entre entités* est un prédicat. Certains parlent de *faits* et d'*événements*.
- Une branche du TAL, le *Semantic Role Labeling* fait l'inverse : étant donné un prédicat, trouver ses différents arguments sémantiques (les réponses aux questions qui-quoi-quand-où-comment, etc.).

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- Les règles peuvent être des simples expressions régulières ou alors des motifs plus complexes, faisant intervenir les mots, les tags POS, des dictionnaires externes, des annotations provenant d'extractions précédentes.
- Une règle typique est du type « motif contextuel → action ».
- L'action peut être le balisage de l'entité, ou d'un groupe d'entités et de relations entre elles.
- L'ordre d'applications des règles est important.
- Les règles peuvent être écrites à la main ou apprises à partir d'un corpus d'apprentissage.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Première approche : le bottom-up.
- [Califf & Mooney 2003] On prend une instance d'entité qui ne soit pas couverte par une règle. À partir de cette instance on crée une règle du type $x_{i-w} \dots x_{i-1} x_i \dots x_{i+w} \rightarrow T$ où T est le tag de x_i . On enlève des tokens ou on les remplace par des tokens plus généraux, et on voit ce que ça donne au niveau du corpus d'apprentissage.
- Deuxième approche : le top-down.
- On part d'une règle très générale et on la spécialise. Par exemple, on part d'une règle R_0 qui pose des conditions au niveau de $2w$ tokens. On va considérer progressivement des règles
 - 1 \mathcal{R}_1 qui diffèrent de R_0 en un des tokens ;
 - 2 pour $L = 2$ à $2w$, des règles \mathcal{R}_L formées par l'intersection de deux règles de \mathcal{R}_{L-1} qui concordent sur $L - 2$ conditions et diffèrent en une.

À chaque fois on regarde la couverture des règles et on ne garde que celles qui dépassent un seuil minimum s .

- [Krishnamurty et al. 2008] définissent un modèle de données relationnel et une algèbre d'opérateurs qui permettent l'optimisation des requêtes. Ils valident leurs résultats sur des corpus de blogs.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- On peut procéder par tokens (mots) ou par segments.
- Dans le premier cas, on part d'une suite de tokens $\vec{x} = x_1 \dots x_n$ et on se propose d'associer à chaque token un label $y_i \in \mathcal{Y}$.
- Dans le deuxième cas, on cherche une suite de segments $s_1 \dots s_p$ et à chaque segment $s_j = l_j \dots u_j$ on associe un label y_i .
- Diverses méthodes ont été employées (HMM, MEMM, CRF). Il paraît que les CRF marchent assez bien.

UVF3B403 MS

IABDA

Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Premier cas : on part de deux entités connues et on cherche leur relations : soit E_1, E_2 entités dans \vec{X} , trouver toutes les relations dans \mathcal{Y} entre E_1 et E_2 .
- On peut écrire des règles (de logique propositionnelle ou de premier ordre) pour détecter les relations.
- On peut utiliser des méthodes à base de propriétés (*features*) :
 - 1 utiliser des listes de mots prédéfinies ou des classes dans des ontologies ;
 - 2 utiliser des tags POS (chercher en particulier : les verbes) ;
 - 3 établir l'arbre syntaxique de la phrase (difficile, et fragile en cas de bruit) ;
 - 4 établir un graphe de dépendances.

Une fois les propriétés choisies, on peut utiliser des arbres de décision, ou des SVM.

- Ou alors, on peut comparer directement des arbres ou des graphes en utilisant des méthodes de noyaux.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Le président de la république₁ a rencontré le serial-killer₂ dans sa_{1 ou 2?} cellule. Il_{1 ou 2?} a avoué éprouver des remords pour ses actes.
- Une **anaphore** est un mot ou un syntagme qui, dans un énoncé, assure une reprise sémantique d'un précédent segment appelé **antécédent**. (La **cataphore** va dans l'autre sens.) On parle aussi de **coréférences**.
- La **résolution d'anaphores** consiste à identifier les entités du texte qui se réfèrent à la même entité du monde réel.
- Il y a différents types de renvois anaphoriques : pronoms personnels (*il, elle,...*), possessifs (*le nôtre,...*), adverbes (*là, ainsi,...*), ordinaux (*le second,...*), etc.
- L'approche générale :
 - ① identifier/délimiter les phrases/paragraphes où l'on va chercher des antécédents ;
 - ② utiliser une série de tests de cohérence (genre, nombre, cas, etc.) pour éliminer un premier lot de candidats ;
 - ③ pondérer, selon certaines règles, les candidats restant ;
 - ④ choisir le candidat de plus grand poids.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Approches à base de règles : CogNIAC [Baldwin 1995] pose six règles, dans un ordre précis. [Kennedy & Boguraev 1996] pondèrent (positivement) selon dix facteurs, impliquant le tag POS et la fonction syntaxique (sujet, coi, cod, etc.). [Mitkov 1998] applique des poids positifs et négatifs selon certains indicateurs (POS, syntaxe, mais aussi distance entre anaphore et antécédent).
- Approches statistiques : [Soon et al. 2001] reprennent l'idée des indicateurs (ils appellent cela des *markables*) et font un apprentissage sur un corpus prétagué.
- Approches sémantiques : selon [Markert & Nissim 2005], alors qu'il y a eu des travaux utilisant un apport sémantique (par exemple : recherche de proximité sémantique dans WordNet, du type *quand mon chien a vu l'os, il l'a mangé*) ceux-ci ne donnent pas de très bons résultats (manque de finesse et de complétude des ontologies).

Exemple : Will Quinlan had not inherited a damaged retinoblastoma supressor gene and, therefore, faced no more risk than other children.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

Réponse à la question : oui, dans Wordnet on a « child < juvenile < person < ... ».

— S: (n) **child**, kid, youngster, minor, shaver, nipper, small fry, tiddler, tike, tyke, fry, nestling (a young person of either sex) « she writes books for children »; « they're just kids »; « 'tiddler' is a British term for youngster »
direct hypernym / inherited hypernym / sister term

— S: (n) juvenile, juvenile person (a young person, not fully developed)

— S: (n) **person**, individual, someone, somebody, mortal, soul (a human being) « there was too much for one person to do »

— S: (n) organism, being (a living thing that has (or can develop) the ability to act or function independently)

— S: (n) living thing, animate thing (a living (or once living) entity)

— S: (n) whole, unit (an assemblage of parts that is regarded as a single entity) « how big is that part compared to the whole? »; « the team is a unit »

— S: (n) object, physical object (a tangible and visible entity; an entity that can cast a shadow) « it was full of rackets, balls and other objects »

— S: (n) physical entity (an entity that has physical existence)

— S: (n) entity (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Après la phase de prétraitement, on a dégagé les informations souhaitées des documents et on les a intégrées dans leurs représentations.
- Il s'agit maintenant de faire du mining sur ces représentations.
- La *classification supervisée* ou *catégorisation* est une classification des documents dans des classes prédéfinies (par exemple : SPAM ou ¬SPAM? sport, culture ou économie?).
- On représente, le plus souvent, les documents par des *vecteurs de propriétés* (features), qui peuvent être pondérés.
- Le modèle du « sac de mots » (*bag-of-words*) consiste à considérer, pour chaque mot w , sa présence (ou sa fréquence) comme une propriété (une dimension de l'espace vectoriel). Dans certains cas, on ajoute aux mots (ou on les remplace par) leurs synonymes / hypéronymes, etc.
- Comme poids du terme t_k dans le document d_j on peut utiliser la fonction *tfidf* :

$$\text{tfidf}(t_k, d_j) = \# \{t_k \in d_j\} \cdot \log \frac{\# \{d_*\}}{\# \{d_* \mid t_k \in d_*\}}$$

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Pb : si chaque mot (ou terme ou concept) du document est une dimension du vecteur, on risque d'en avoir *beaucoup*. Différentes méthodes ont été appliquées pour réduire le nombre de dimensions.
- Une méthode linéaire est le LSI (*latent semantic indexing*) [Deerwester et al. 1990] :
- un résultat d'algèbre linéaire dit que :
 - 1 toute matrice X de rang n peut s'écrire USV^t , où U et V sont orthogonales ($UU^t = I$, $VV^t = I$) et S est diagonale, de termes diagonaux $s_1 \geq s_2 \geq \dots \geq s_n$;
 - 2 soit $k < n$, en prenant S' diagonale avec les k premiers termes de S , U' les k premières colonnes de U , V' les k premières colonnes de V , on obtient $X' = U'S'V'$;
 - 3 alors X' est la matrice de rang k la plus proche (norme de Frobenius) de X .
- On applique cette méthode (appelée « décomposition en valeurs singulières ») à la matrice S des termes \times documents, et on obtient S' de rang k , optimale parmi toutes les matrices de ce rang.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- LSI est une méthode *linéaire* de réduction de la dimensionnalité.
- Autres approches, provenant du traitement de l'image : [Tenenbaum et al. 2001] utilisent, au lieu de la distance euclidienne, des distances géodésiques sur des variétés topologiques qui collent au mieux aux données (*manifold learning*).
- L'approximation de ces variétés se fait à l'aide de graphes obtenus en liant par des arêtes des points de l'échantillon avec les voisins se trouvant à distance (euclidienne) inférieure à ε . Ensuite, une fois le graphe obtenu, on considère comme distance entre deux points, la longueur du plus court chemin (Dijkstra). C'est l'algorithme *ISOMAP*.
- Après, on réduit la dimension du graphe de manière classique.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- Une fois les documents représentés par des vecteurs de dimension raisonnable, on utilise les méthodes classiques de data mining.
- Approche par règles logiques (DNF) écrites manuellement [Hayes et al. 1990].
- Approches statistiques : le naïve Bayes est assez robuste, malgré l'hypothèse que les mots soient tous indépendants [Domingos & Pazzani 1997].
- Regression logistique, arbres de décision (ID3, C4.5, CART), règles DNF apprises sur un corpus, méthode de Rocchio, réseaux de neurones, kNN, SVM, AdaBoost, etc.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- La *classification non supervisée* ou *clustering*, consiste à classer les documents sans connaissance préalable des catégories.
- Il s'agit de comparer toutes les partitions possibles d'un ensemble de documents à l'aide d'une *mesure de similarité* : les documents dans le même groupe doivent être aussi similaires que possible, les groupes entre eux aussi dissimilaires que possible.
- Première approche : *Hierarchical Agglomerative Clustering* :
 - ① on met chaque objet dans un cluster séparé ;
 - ② on fusionne les clusters les plus similaires ;
 - ③ on répète le (2) jusqu'à ce que tout soit dans le même cluster.
- [El-Hamdouchi & Willet 1989] Il y a différentes manières de mesurer la similarité des clusters : *single-linkage* (max de similarité entre leurs obj), *complete-linkage* (min), *centre de gravité* (similarité entre centroïdes), *average group* (similarité moyenne), *Ward* (plus petite augmentation de distance du centroïde du nouveau cluster), etc.
- Le pb est que l'on ne peut pas revenir en arrière pour corriger

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- Deuxième approche : l'algorithme *K-means* :
 - 1 on prend au hasard k vecteurs et on considère qu'ils sont les centres des clusters : tout vecteur appartient au cluster dont le centre est le + proche ;
 - 2 on calcule les *centroïdes* : les moyennes des tous les vecteurs de cluster, ils deviennent les n centres de cluster et on répartit de nouveau les vecteurs ;
 - 3 on répète le (2) jusqu'à convergence.
- Le K-means dépend bcp du choix initial.
- À l'aide du LSI (latent semantic indexing) on réduit le nombre de dimensions et on obtient, non pas des centroïdes, mais des *médoïdes* : des documents artificiels qui jouent le même rôle.

UVF3B403 MS

IABDA

Fouille de texte

Yannis

Haralambous
(IMT Atlantique)

ABox

TBox

RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- [Hotho *et al.* 2003] décrivent un clustering enrichi par des accès à des ontologies ou à Wordnet. Ils décrivent plusieurs stratégies :
 - 1 « add » à tout vecteur de termes (un vecteur représentant le document, dont les coordonnées correspondent à la fréquence des termes) on ajoute un vecteur de concepts pour ceux qui ont pu être identifiés,
 - 2 « repl » quand un terme renvoie à un concept, on remplace les coordonnées de termes par celles des concepts,
 - 3 « only » on n'utilise que des concepts.
- La correspondance entre termes et concepts étant ambiguë, ils donnent également des stratégies de désambiguïsation :
 - 1 « all » on prend tous les concepts liés au terme,
 - 2 « first » si on a un ordre d'importance des concepts, on ne prend que le premier,
 - 3 « context » on désambiguïse par la méthode de la « densité de concepts » : on prend les sur- et sous-concepts et les termes qui les représentent et on regarde leurs fréquences dans le document.

- Première approche : par la *compréhension* ;
- inspirée des sciences cognitives et de l'IA ;
- on construit une représentation sémantique du texte, on la réduit et à partir de la réduction on génère un résumé ;
- pour la réduction on utilise des *marco-règles* :
 - 1 l'*élimination* : Pierre a vu une balle bleue → Pierre a vu une belle.
La balle était bleue → Pierre a vu une balle,
 - 2 la *généralisation* : Pierre a vu un faucon → Pierre a vu un oiseau ;
Pierre a vu un faucon. Pierre a vu un vautour → Pierre a vu des oiseaux,
 - 3 la *condensation* : Pierre a creusé le fondations, construit les murs, posé le toit... → Pierre a construit une maison ;
- problèmes : construire la représentation sémantique peut être couteux ; choisir ce qui est important lors de la phase de réduction peut être délicat ; on est au sommet de l'édifice, il ne faut pas se planter dans les couches inférieures.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- Deuxième approche : par *extraction* ;
- on trouve les phrases les plus importantes et on les extrait.
- Première idée : trouver les phrases dont les mots on le td-idf global le plus élevé (suppose l'existence d'un corpus). Facile à calculer, mais est-ce bien les phrases les plus intéressantes ?
- problèmes : pour des raisons stylistiques, les auteurs cultivent la synonymie ; les anaphores ne sont pas résolues ; il faut au moins une lemmatisation pour que les calculs soient cohérents ; quid de l'ordre des phrases ?
- Deuxième idée : se baser sur des phrases prototypiques : *dans cet article nous allons nous intéresser à ...* en affectant un score selon la présence ou non de certains mots indicateurs ;
- problèmes : demande une adaptation à chaque type de texte ; les anaphores ne sont pas résolues ; cohérence du résultat ?

UVF3B403 MS
IABDA
Fouille de texteYannis
Haralambous
(IMT Atlantique)ABox
TBox
RBoxConstructeurs et
restrictionsConstructeurs de
concepts
Restrictions de rôleOntologies DL
SROIQDéfinitions et
exemples
Requêtes

- Troisième idée : par repérage de chaînes lexicales ;
- on considère, un par un, les noms. On cherche leurs distances dans WordNet et on cherche un graphe de relations entre eux qui corresponde au mieux à ces distances : il y a des fortes chances à ce que ce soient les relations qu'ils ont dans le texte ;
- ensuite, à l'aide de ces relations on attribue un score aux phrases et on prend celles de plus grand score.
- Quatrième idée : en analysant les relations entre les phrases ;
- après des analyses morphologique et syntaxique, on cherche certains patrons (par exemple, alors, par conséquent, etc.) qui permettent d'établir le graphe des relations entre phrases, et puis entre paragraphes ;
- ensuite, dans ce graphe, on trouve les nœuds (phrases) qui ont des « rôles sémantiques » précis, par exemple celui de conclusion.

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts

Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples

Requêtes

- Cinquième idée : en construisant un « analyseur rhétorique »;
- [Marcu, 1997] a construit une base de 450 marqueurs discursifs qui permettent d'établir des relations discursives entre les propositions d'un texte ;
- il a développé un algorithme qui construit un arbre optimum dont les flèches correspondent à des relations rhétoriques (du type élaboration, justification, exemplification, concession, antithèse, contraste, évidence, etc.) ;
- la racine de l'arbre sera la proposition la plus saillante et on fait un parcours en largeur d'abord jusqu'à atteindre la taille de résumé souhaitée.

- Troisième approche : par *apprentissage* ;
- dans le cas de l'apprentissage supervisé, on prend un corpus, on en extrait des phrases (ou on leur attribue un score d'extraction) et on écrit des critères d'extraction. Ceux-ci peuvent être
 - ➊ positionnels (la première phrase est plus importante, etc.),
 - ➋ morphologiques et quantitatifs (comme les fréquences de certains termes),
 - ➌ discursifs ;
- à partir des critères d'extraction, on construit pour chaque phrase du corpus d'entraînement un vecteur de valeurs de critères ;
- l'algorithme compare ces vecteurs avec le score d'extraction donné manuellement, pondère les règles, et évalue les résultats sur un corpus de test ;
- en introduisant des nouvelles règles, on peut améliorer les performances du système ;
- on peut même imaginer des apprentissages non supervisés ou semi-supervisés.
- Problème : sans analyse morphosyntaxique ni résolution d'anaphores ; l'instance de l'apprentissage est la phrase, les relations entre phrases

UVF3B403 MS
IABDA
Fouille de texte

Yannis
Haralambous
(IMT Atlantique)

ABox
TBox
RBox

Constructeurs et
restrictions

Constructeurs de
concepts
Restrictions de rôle

Ontologies DL
SROIQ

Définitions et
exemples
Requêtes

- À travers quelques notions de base et quelques exemples, nous avons fait le tour des enjeux et méthodes de base de la fouille de texte ;
- mis à part la complexité du phénomène « langue », ce qui rend le sujet encore plus vaste est le fait qu'il n'y a pas de véritable frontière entre les différentes disciplines : linguistique, traitement automatique de langue, apprentissage artificiel, statistique, algorithmique, etc.
- ce qu'il faut retenir de ce cours : les différents angles d'attaque, les avantages et inconvénients des différentes méthodes et la tendance générale de les essayer toutes, en les comparant par des techniques d'évaluation ;
- plus spécifiquement, concernant la langue, les deux approches complémentaires : celle des méthodes formelles (logique, λ -calcul, etc.) et celle des méthodes statistiques, souvent avec apprentissage. Les meilleurs résultats s'obtiennent en les combinant.