



UNIVERSITÉ  
DE MONTPELLIER



# Soutenance TER M1 : Hotel Advisor

Outils d'analyse automatique d'avis d'hôtel

---

## Étudiants

Bastien	CARBONNIER
Cedric	DURAND
Johann	GOLMARD

## Encadrants

Mathieu	LAFOURCADE
Pierre	POMPIDOR

4 juin 2019

## Sujet

Réaliser une application qui analyse automatiquement des avis textuels sur des hôtels pour ensuite « colorer » les différents descripteurs des ontologies représentant ces différents hôtels.

## Avis

la chambre était très sale

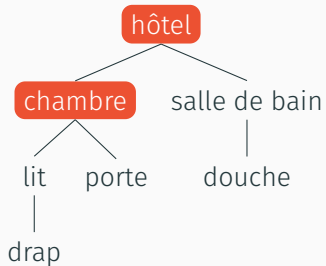
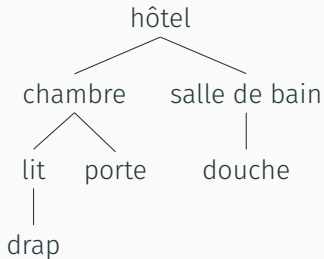
## Après polarisation

la chambre était très sale

## Après propagation

la chambre était très sale

## Ontologies :



1. État de l'art
2. Outils utilisés
3. Architecture de notre projet
4. Limites et évolutions

## État de l'art

---

## Définition

Consiste à identifier pour chaque mot sa nature grammaticale à partir du contexte de la phrase dans laquelle il se situe et de connaissances lexicales.

- Un mot peut avoir plusieurs natures
- Bibliothèques Python TreeTagger, SpaCy, NLTK
- **Important** : Base de notre propagation de polarisation

## Protocole

- 15 Avis réels (2331 mots)
- Redressement orthographique
- Étiquetage grammatical en utilisant les 2 librairies
- 291 termes étiquetés différemment
- Recherche des bonnes natures pour ces termes
- Comparaison et interprétation des résultats

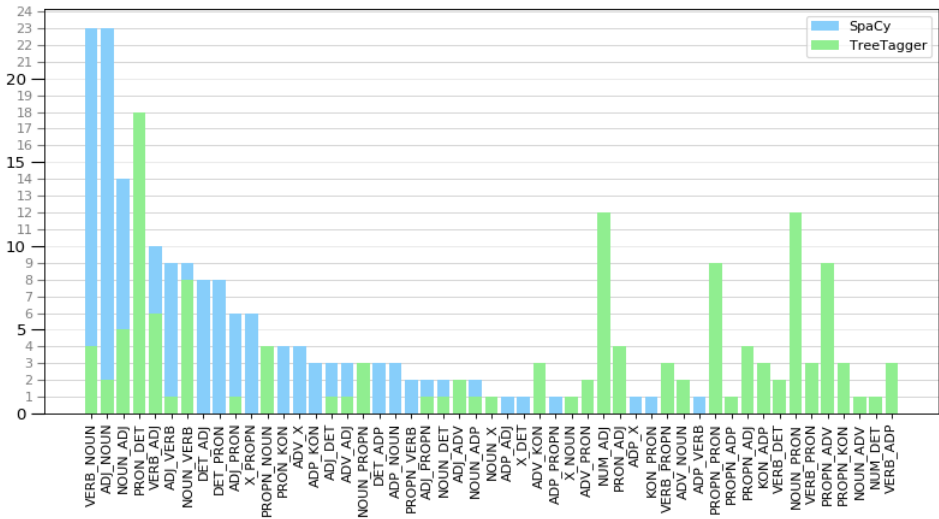
## Précisions :

- TreeTagger =  $\simeq$  53%
- SpaCy =  $\simeq$  46%



# Comparaison TreeTagger et SpaCy

Comparatif erreur TreeTagger SpaCy (tagFound\_tagRef)



# Comparaison TreeTagger et SpaCy

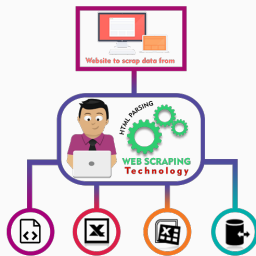
SpaCy (pred_ref)			TreeTagger (pred_ref)		
verb_noun	adj_noun	noun_adj	pron_det	noun_adj	noun_pron
mer	jardins	junior	chaque	futée	Il
déjeuner	safari	grossier	cet	troisième	tout
coucher	buffets	bel	les	solo	J
bar	armoire	sympa	cette	petits	Tous
déjeuner	pente	chaude	le	moderne	Ils
forfait	bars	superbe	les		
fruits	magasins	gratuite			
jacuzzis	restaurants	cool			
vigile	vagues	correct			
serviettes	lits	médiocre			
Suites	bouilloire	abordables			

**Table 1** – Mots mal étiquetés

## Conclusion

- Erreurs d'étiquetages de TreeTagger plus facilement réparables
- Erreurs d'étiquetages de SpaCy impactent la sémantique de la phrase
- Étude relative et non complète

# Scraping



```
135 4.0|Ce que l'on a aimé : Le all inclusive avec un grand choix de plats et des alcools internationaux. La proximité du
contre commercial Yumbo. Très grande chambre studio avec belle terrasse et vue mer, mais un peu surannée, un solarium
nautiste sur le toit de l'hôtel.
136 4.0|Bien dans l'ensemble, nous avons passé une semaine agréable, la piscine est assez petite quand aux transats il faut se
lever de bonheur pour pouvoir en bénéficier car il sont prix d'assaut.A revoir l'insonorisation des chambres qui fait
défaut
137 4.0|hôtel impeccable, jardins magnifiques et entretenus tous les jours, nourriture variée et excellente(notamment les jus
de fruits pressés le matin)l'iterie très confortable.petit bébé, peu de Français, donc, peu de personnel le parlant...a
conseiller, la formule demi pension suffit largement à mon avis
138 4.0|Avons passé une semainebuffet varié Personnel agréable Animation présente , de beaux spectacles le soir.Dommage que le
programme n'est pas diffusé sur papier dans les chambres.Beaucoup de monde mais 2 piscines dont 1 au calme
139 4.0| bon rapport qualité prix très peu de chose à dire sur l'ensemble des prestations,nous sommes satisfaits aucune
remarque, peut être l'animation à revoir qui pour moi est leur point faible,nous étions 3 adultes et 3 enfants dont 2
bébés tout était parfait pour ce séjour d'une semaine en famille
140 4.0|J'ai voyagé à travers toute la planète et les îles canaries sont parmi les plus belles destinations. Le personnel de
occidental margaritas est formidable. Un professionnalisme hors pair, une incroyable gentillesse. Parmi les meilleurs
services au monde. Malheureusement la clientèle n'est pas au même niveau. Je me suis fait agresser par un espagnol dont la
mère a prétendu qu'il était policier pour m'impressionner. J'ai réclmé à la réception et j'espère qu'ils ne vont pas
prendre le parti de leur compatriote. Ce serait vraiment dommage. En tout cas ces messieurs dames devraient savoir qu'on
ne peut être impressionné par la police espagnole que quand on a quelque chose à se reprocher. Dans leur cas ce devrait
être la mauvaise éducation. Le polo que ce messieur à jeté à la piscine vaut beaucoup plus que ses manières de bas fond
141 4.0|hotel joli un peu en travaux en ce moment piscine bien restaurant pas mal mais ça fais trop usine! vu que cest all
inclus personnel sympa tres bien situé pour les gays cest proche de yumbo 5 min et 10 min de la plage donc cest
excellent!l'accueil sympa des chats à l hotel trop mignon hihh
142 5.0|excellent établissement dépayçant, cadre idyllique, personnel dévoué, et à l'écoute de nos souhaits, aimables et
souriants, agréables sympathique. Établissement propre, et toujours bien entretenu.Très bon souvenirs, y reviendrais bien
volontiers.
```

## Définition

- Outils de représentation d'un corpus
- Ensemble structuré de concepts sous forme de graphe
- Relations sémantiques et/ou de composition et héritage

## Termes de l'ontologie

- Domaine restreint à l'hôtellerie
- Chaque terme colorié selon sa polarité
- Ontologie fournie par notre encadrant après formatage

## Format JSON

- Lisibilité aisée
- Représentation arborescente
- Compatible avec MongoDB



## Représentation avec D3.js

- Permet de lier des données à un objet DOM
- Offre de nombreuses interactions personnalisables et fluides
- Librairie présente sur Angular, qui peut être découpée en sous-modules
- Documentation complète et conséquente



## Outils utilisés

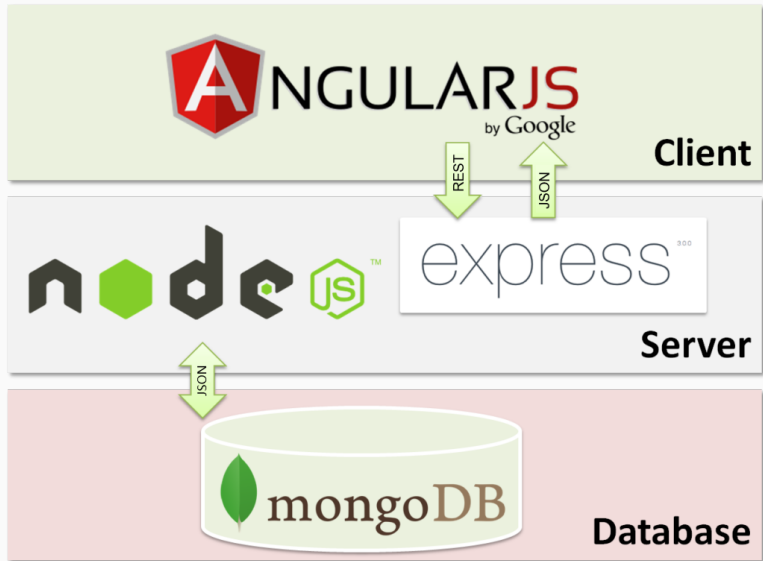
---



- Réseau lexical
- GWAP (Game with a purpose)
- Point d'accès rezo\_dump.php
- Utilisation d'un cache



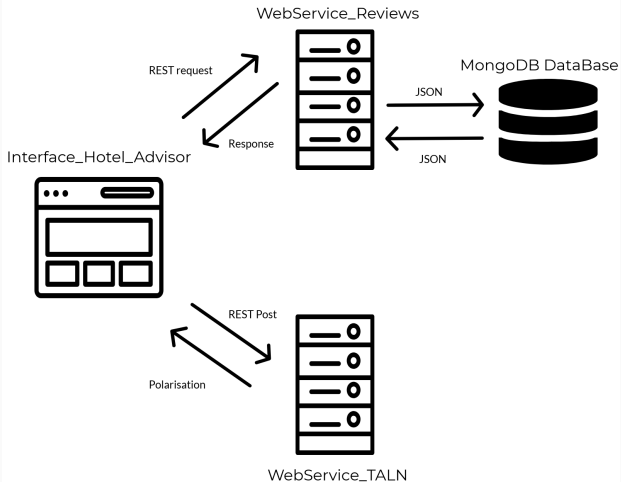
# Architecture MEAN



## Architecture de notre projet

---

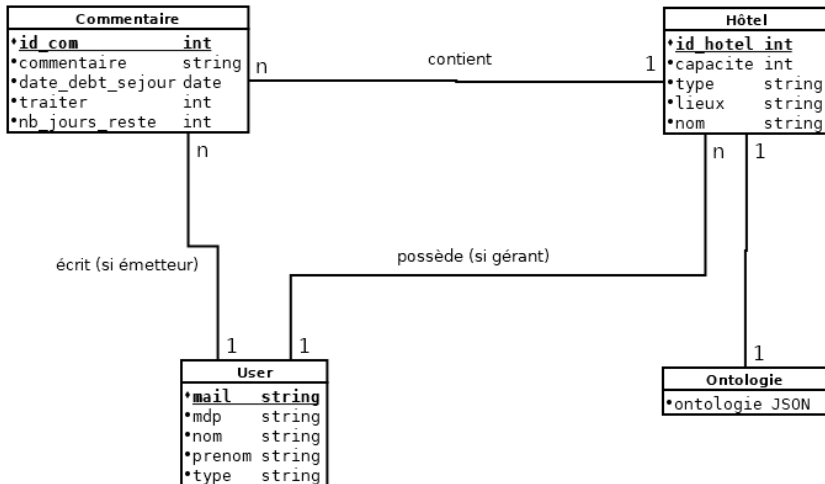
# Schéma global



## Webservice BDD

---

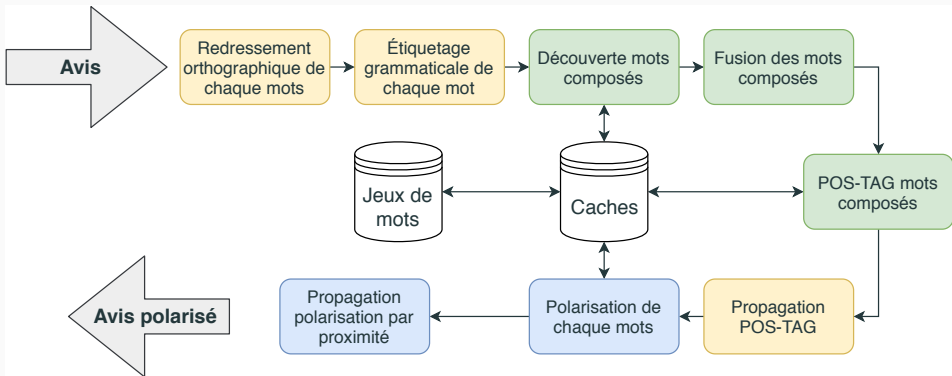
# Webservice BDD



## Webservice TALN

---

# Démarche adoptée





## Correcteur orthographique

- Module Python spellchecker
- Basé sur la distance de Levenshtein
- Recherche dans une liste contenant la fréquence probable d'un terme

## Lemmatisation

- Effectuée par la librairie d'étiquetage grammaticale
- Utilisation du lemme si le mot n'existe pas dans JeuxDeMots

# Étiquetage morphosyntaxique

## Librairie choisie : TreeTagger

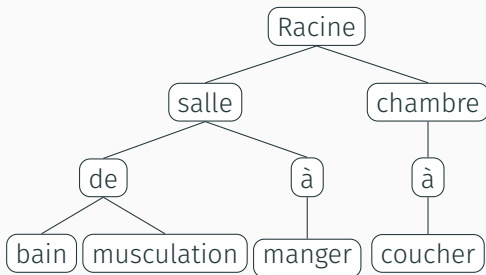
- Meilleur étiquetage grammaticale
- Plus rapide : SpaCy ( $\simeq 10s$  par avis), TreeTagger ( $\simeq 4s$  par avis)

## Structure de données utilisée

```
1 [
2   {index : 0, mot : "la", nature : "DET", lem : "le"},
3   {index : 1, mot : "chambre", nature : "NOUN", lem : "chambre"},
4   {index : 2, mot : "était", nature : "VERB", lem : "être"},
5   {index : 3, mot : "sale", nature : "ADJ", lem : "sale"}
6 ]
```

## Protocole

- Création d'un arbre des préfixes des mots composés
- Recherche pour chaque position  $i$  dans la phrase le mot composé le plus grand ayant pour premier mot `phrase[i]`
- Fusion des mots composés
- Recherche des étiquettes grammaticales de ces nouveaux mots sur JDM



# Propagation des étiquettes grammaticales

## Protocole

- Pour chaque NOM recherche des ADJs et du VERB
- Pour chaque ADJ recherche des ADVs
- Pour chaque VERB recherche des ADVs

## Exemple

```
1  [  
2    {index : 0, mot : "la", nature : "DET"},  
3    {index : 1, mot : "chambre", nature : "NOUN",  
4      index_verbe : [2],  
5      index_adj : [4]},  
6    {index : 2, mot : "était", nature : "VERB"},  
7    {index : 3, mot : "très", nature : "ADV"},  
8    {index : 4, mot : "sale", nature : "ADJ",  
9      index_adv : [3]}  
10 ]
```

## Idées

1. Gérer la polarité comme un entier (-1, 0 ou +1)
2. Gérer la polarité comme un vecteur

## Exemple

```
1  [ {index : 0, mot : "la", nature : "DET",  
2    pol: { neg: 0, pos: 0, neutre: 0 }},  
3    {index : 1, mot : "chambre", nature : "NOUN",  
4      index_verbe : [2], index_adj : [4],  
5      pol: { pos: 0.16, neutre: 0.84, neg: 0 }},  
6    {index : 2, mot : "était", nature : "VERB",  
7      pol: { pos: 0.75, neutre: 0.16, neg: 0.09 }},  
8    {index : 3, mot : "très", nature : "ADV"},  
9    {index : 4, mot : "sale", nature : "ADJ",  
10     index_adv : [3],  
11     pol: { pos: 0.02, neutre: 0.02, neg: 0.96 }}}
```

# Propagation de la polarisation

## Ordre de propagation

- $ADV \rightarrow ADJ, ADV \rightarrow VERB$  (Moyenne)
- $ADJ \rightarrow NOM, VERB \rightarrow NOM$  (Moyenne ou substitution)

## Exemple

```
1  [ {index : 0, mot : "la", nature : "DET",
2    pol: { neg: 0, pos: 0, neutre: 0 }},
3    {index : 1, mot : "chambre", nature : "NOUN",
4      index_verbe : [2], index_adj : [4],
5      pol: { pos: 0.01, neutre: 0.01, neg: 0.98 }},
6    {index : 2, mot : "était", nature : "VERB",
7      pol: { pos: 0.9, neutre: 0.06, neg: 0.04 }},
8    {index : 3, mot : "très", nature : "ADV"},
9    {index : 4, mot : "sale", nature : "ADJ",
10     index_adv : [3],
11     pol: { pos: 0.01, neutre: 0.01, neg: 0.98 }}}
```

## Protocole

- 14 avis
- 82 termes de notre ontologie

## Résultats

- Accuracy = 51,12
- Précision<sub>pos</sub> = 0.67
- Précision<sub>neg</sub> = 1
- Rappel<sub>pos</sub> = 0.69
- Rappel<sub>neg</sub> = 0.22

		Références		
		1	0	-1
Dédit	1	34	0	17
	0	15	1	8
	-1	0	0	7

**Table 2** – Matrice de confusion

Salle de bain sale, baignoire bouchée et poils dans la grille... Cuvette des WC jaunit par l'urine et les années... Planché bancal et bruyant... Aucune intimité cloison non isolées et un jour de 10 cm entre le sol et la porte d'entrée... Sans parler de la prestation restaurant : accueil et service très moyen, plats sans goûts et pourtant la maison vante les produits du terroir... (indigestion garantie...) et pour terminer un chien couché à l'entrée qui empest l'accueil et nous empêche le passage au comptoir.

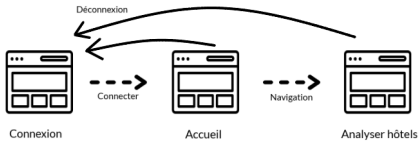


# Application Angular

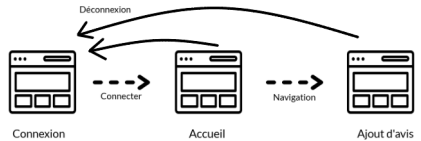
---

# Workflow

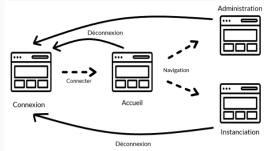
## Gérant



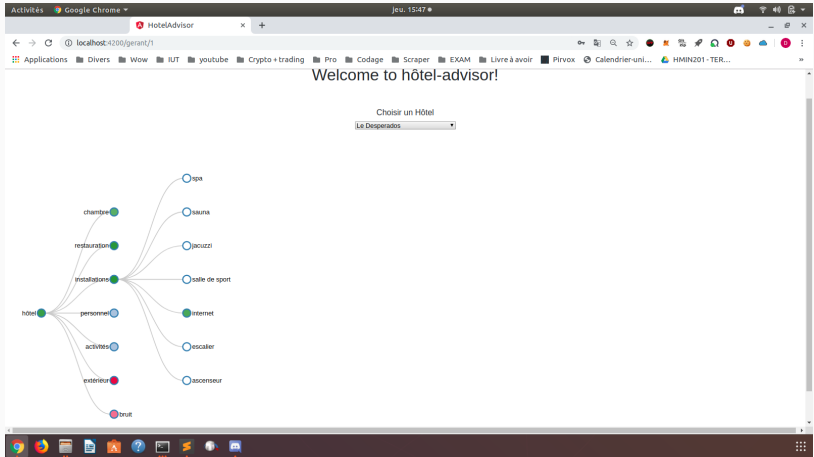
## Émetteur



## Admin



# Ontologie



# Limites et évolutions

---

## Limites

- Fautes de syntaxe très difficilement réparables
- Fautes d'orthographe
- Étiquetage grammaticale imparfait

## Évolutions

- Amélioration de l'interface
- Amélioration des pré-traitements
- Utilisation d'une autre propagation de polarité en utilisant un arbre syntaxique
- Gestion temporelle de l'ontologie