



Réseau IN-OVIVE

Approche « ontologico-vectorielle » pour l'extraction de termes d'intérêts dans la littérature scientifique

Arnaud Ferré

MaIAGE, INRA, Université Paris-Saclay
LIMSI, CNRS, Université Paris-Saclay

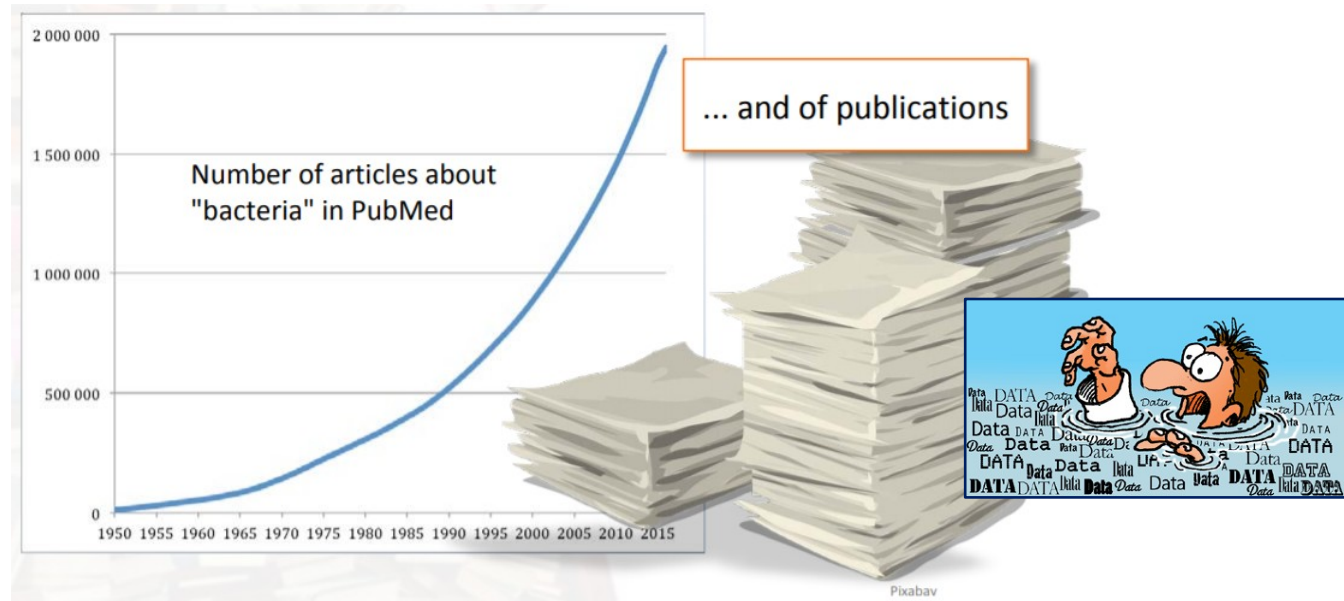


INTRODUCTION

Problématique générale :

Beaucoup de connaissances sont sous forme non-structurée dans des textes (publications scientifiques, champs de texte dans une BDD, sites web, etc.)

Exemple :



INTRODUCTION

Définition :

L'**extraction d'information** vise à analyser du texte (*i.e.* **non-structuré**) pour interpréter et construire une représentation **formelle** d'une partie de l'information contenu dans celui-ci.

Les représentations formelles générées permettent de faire le lien avec d'autres sources de données formalisées.



INTRODUCTION

Reconnaissance
d'entités nommées

Normalisation

Extraction de
relations

[Bacteria]

[Bacteria]

M. agassizii and *M. Testudineum* are present in

Georgia populations of gopher tortoises.

[Geographical]

[Habitat]

INTRODUCTION

Reconnaissance
d'entités nommées

Normalisation

Extraction de
relations

[Bacteria]

NCBI 33922

[Bacteria]

NCBI 244584

M. agassizii and *M. Testudineum* are present in

Georgia populations of gopher tortoises.

[Geographical]

Georgia (GE)

[Habitat]

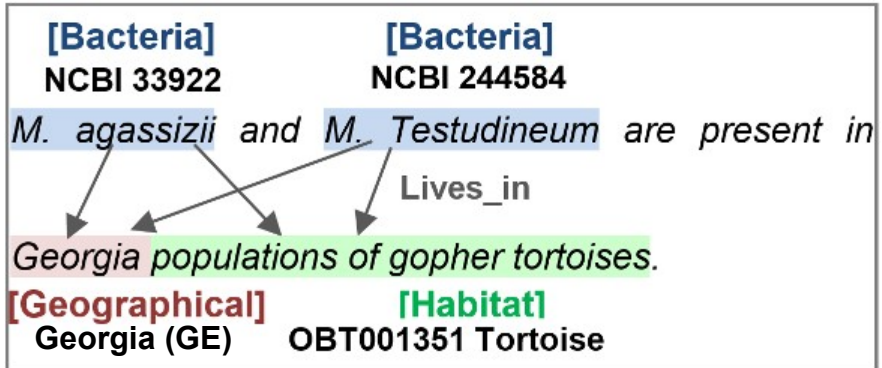
OBT001351 Tortoise

INTRODUCTION

Reconnaissance
d'entités nommées

Normalisation

Extraction de
relations

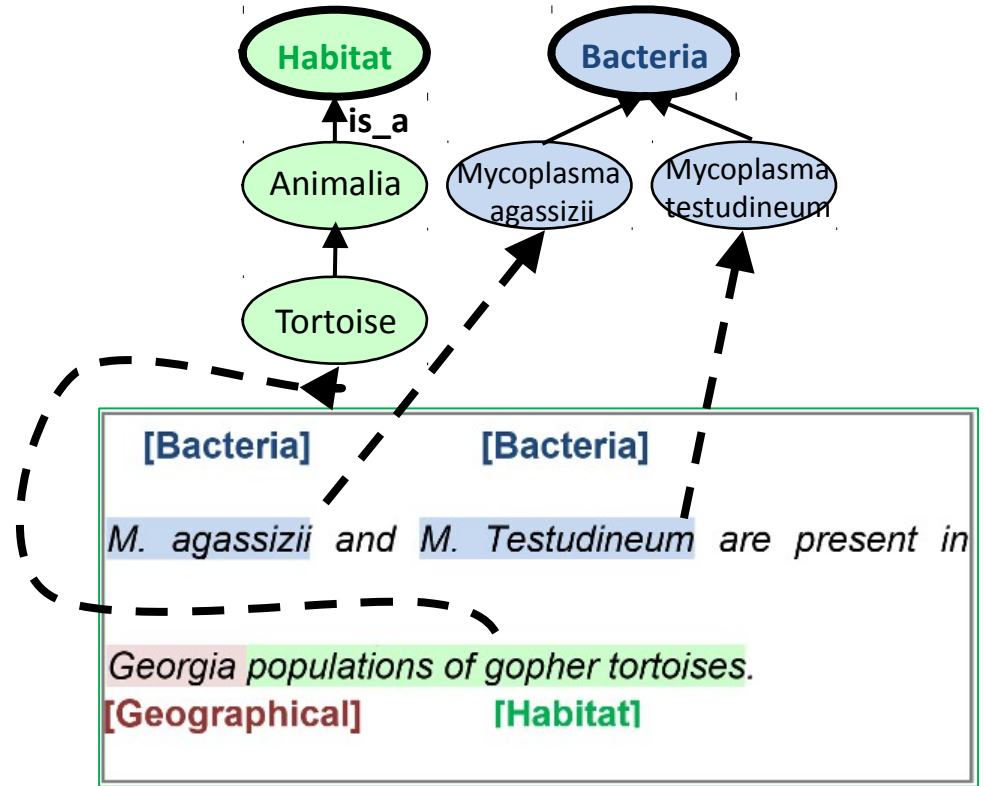


INTRODUCTION

Normalisation

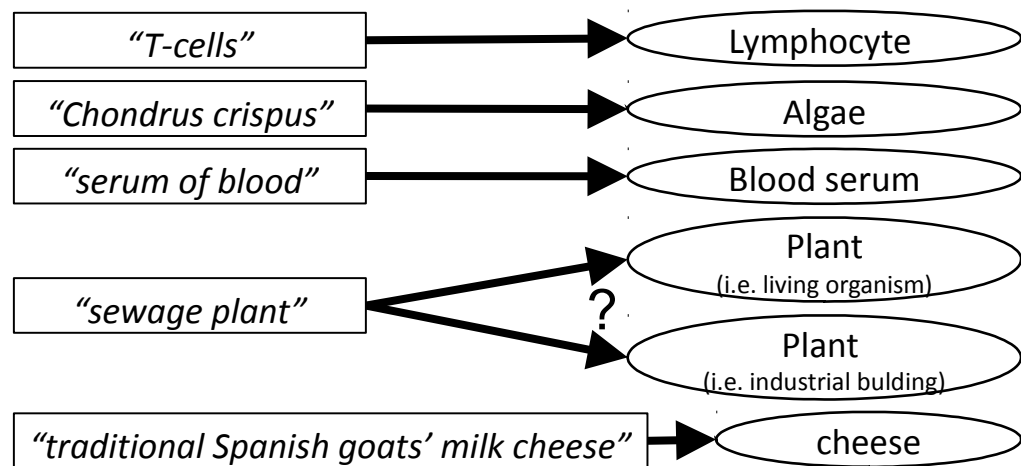
Méthode :

- par des concepts d'une ontologie
- applicable à tout domaine



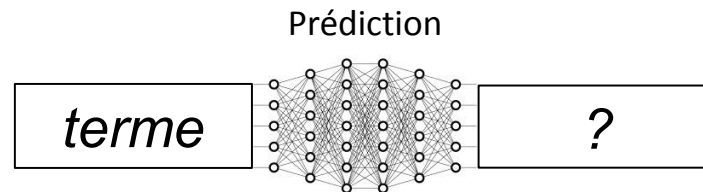
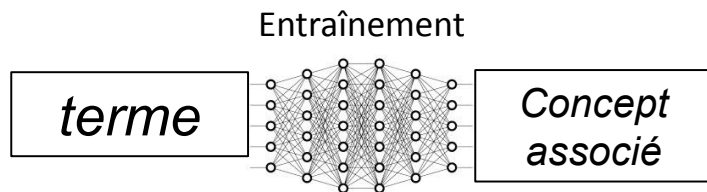
DIFFICULTÉS ET ETAT DE L'ART

- Forme et structure variables
- Approche par similarité de forme (*rule-based, dictionary-based, ...*) ❌
(précision+, rappel-, domaine-spécifique)



- Beaucoup de concepts-cibles, données annotées rares, couverture faible (i.e. peu de concepts utilisés), mots avec faible fréquence d'apparition dans les corpus
 - OntoBiotope : > 2000 concepts
 - UMLS : > 3 millions !

- Approche par classification "directe" supervisée ❌



MATÉRIEL

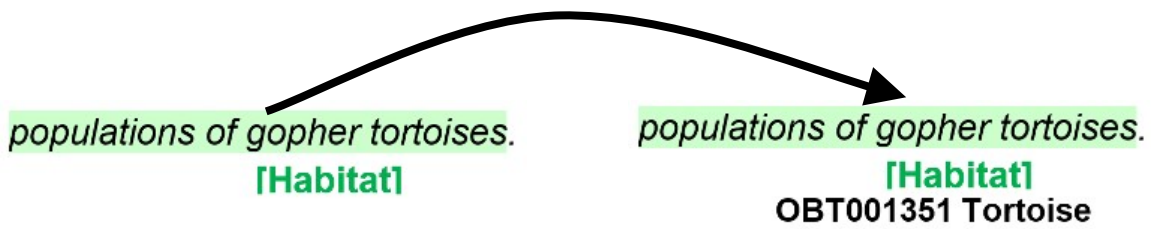
populations of gopher tortoises.

[Habitat]

populations of gopher tortoises.

[Habitat]

OBT001351 Tortoise



Description :

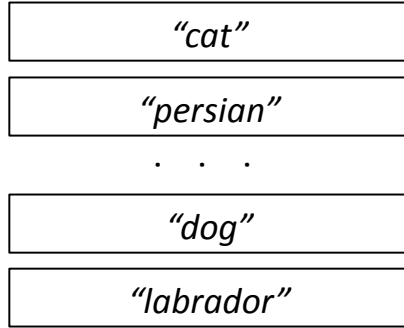
- Challenge Bacteria Biotope de BioNLP Shared Task 2016
- Objectif : Normalisation de mentions d'habitat bactérien
- Données : Sur des titres + résumés d'articles scientifiques en biologie

3 Corpus :

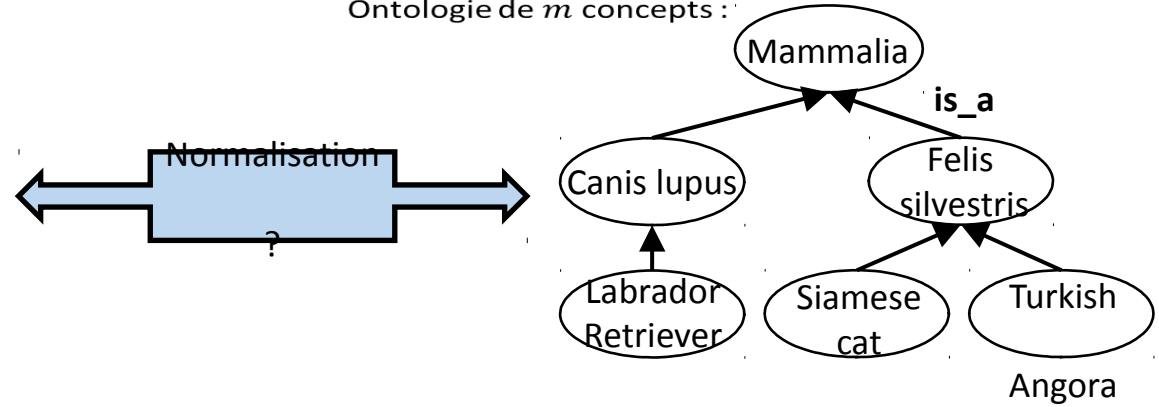
- Corpus d'entraînement : termes d'habitat bactérien + concepts associés (< 1500 associations annotées manuellement)
- Corpus élargi non-annoté du domaine biomédical (PubMed)
- Corpus de test : juste les termes d'habitat bactérien

MÉTHODE : Représentation vectorielle One-Hot

Vocabulaire de n termes :

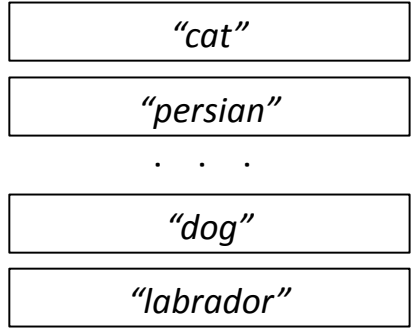


Ontologie de m concepts :

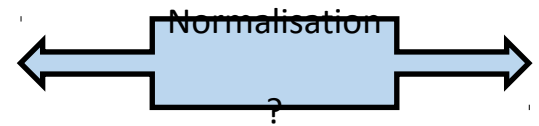
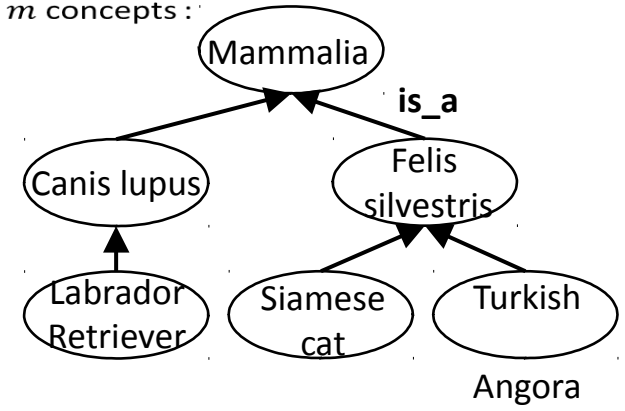


MÉTHODE : Représentation vectorielle One-Hot

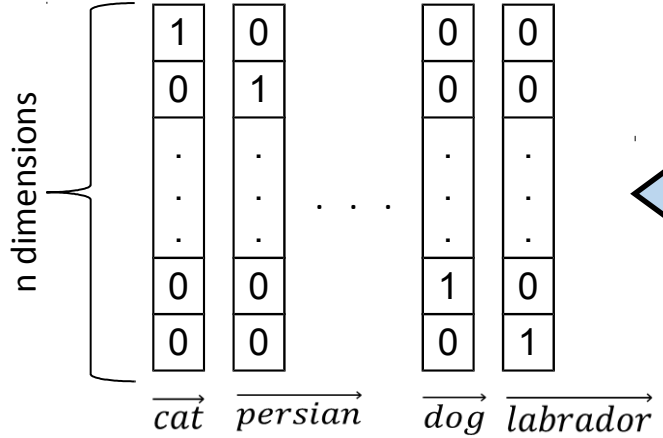
Vocabulaire de n termes :



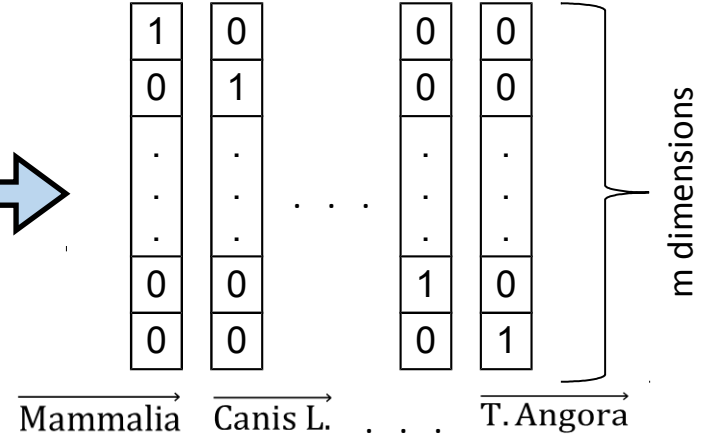
Ontologie de m concepts :



One-hot pour mots :

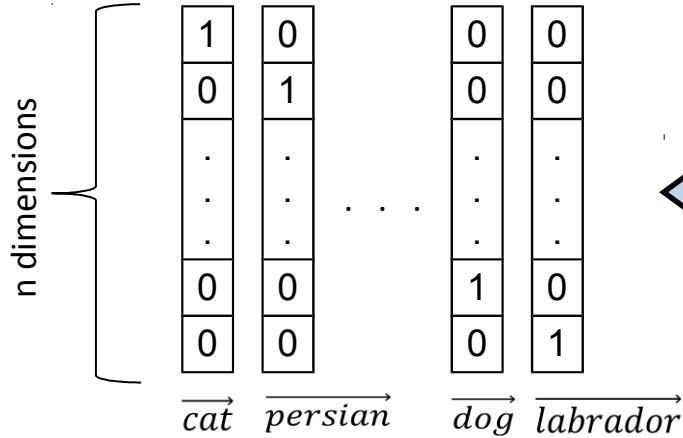


One-hot pour concepts :

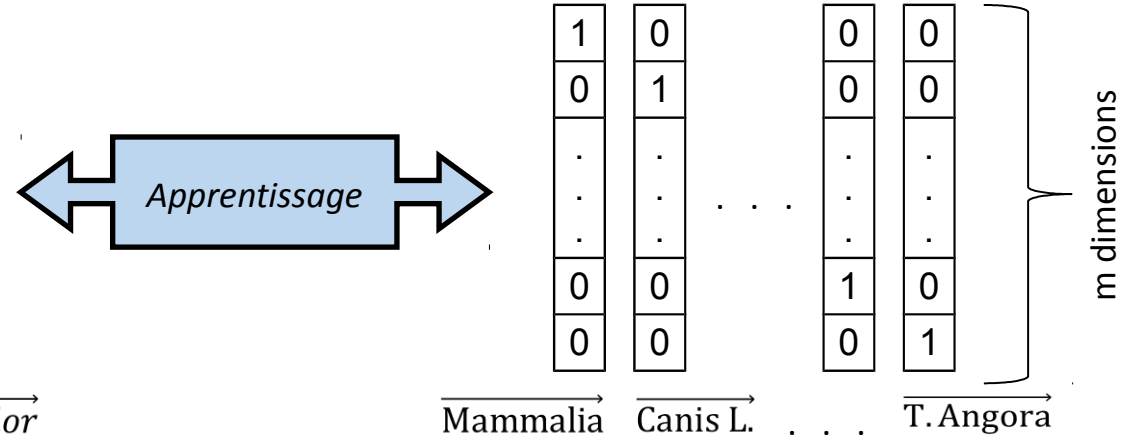


MÉTHODE : Représentation vectorielle One-Hot

One-hot pour mots :



One-hot pour concepts :



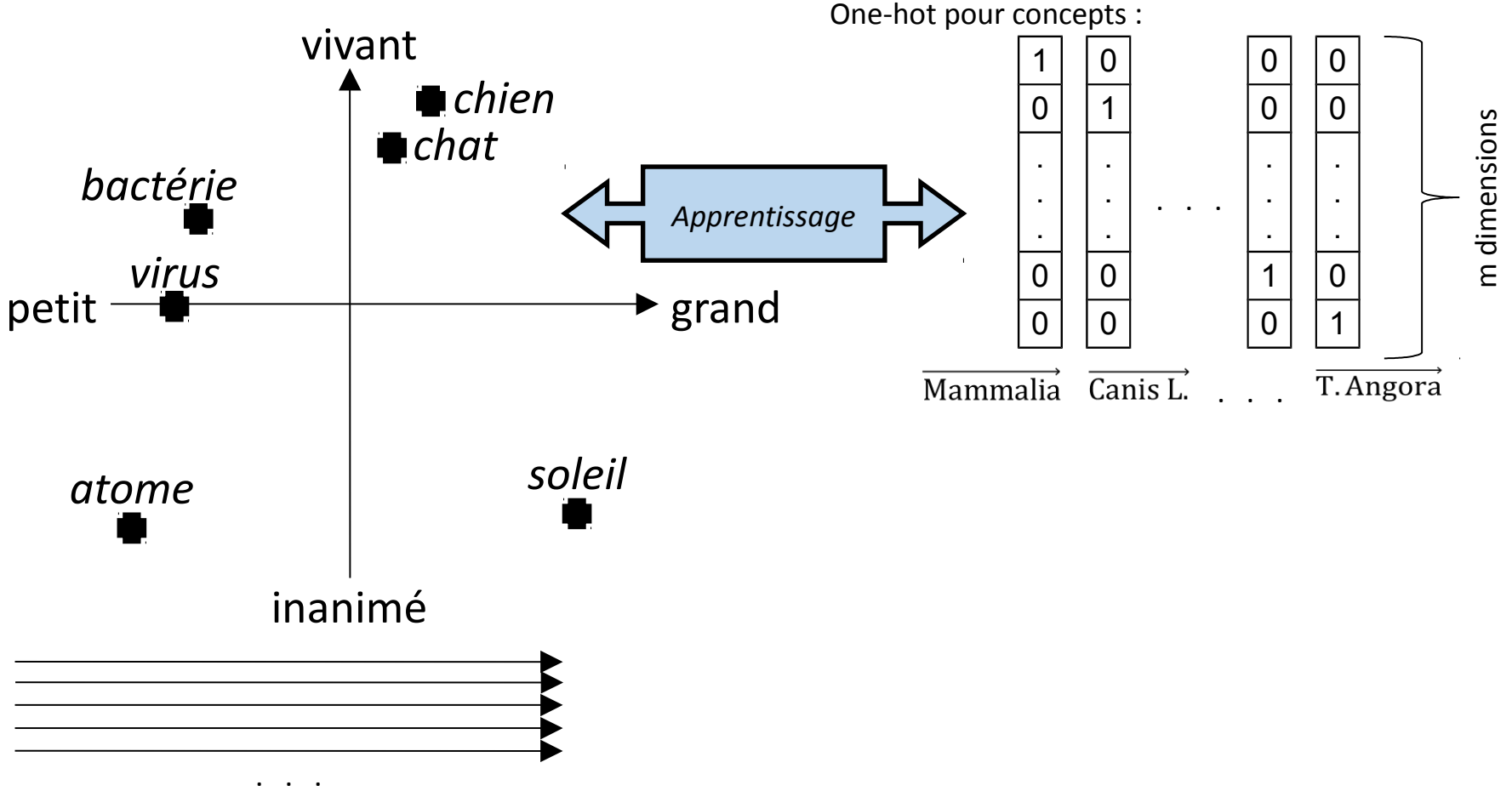
Avantage :

- Représentation vectorielle = nécessaire pour la plupart des algorithmes d'apprentissages

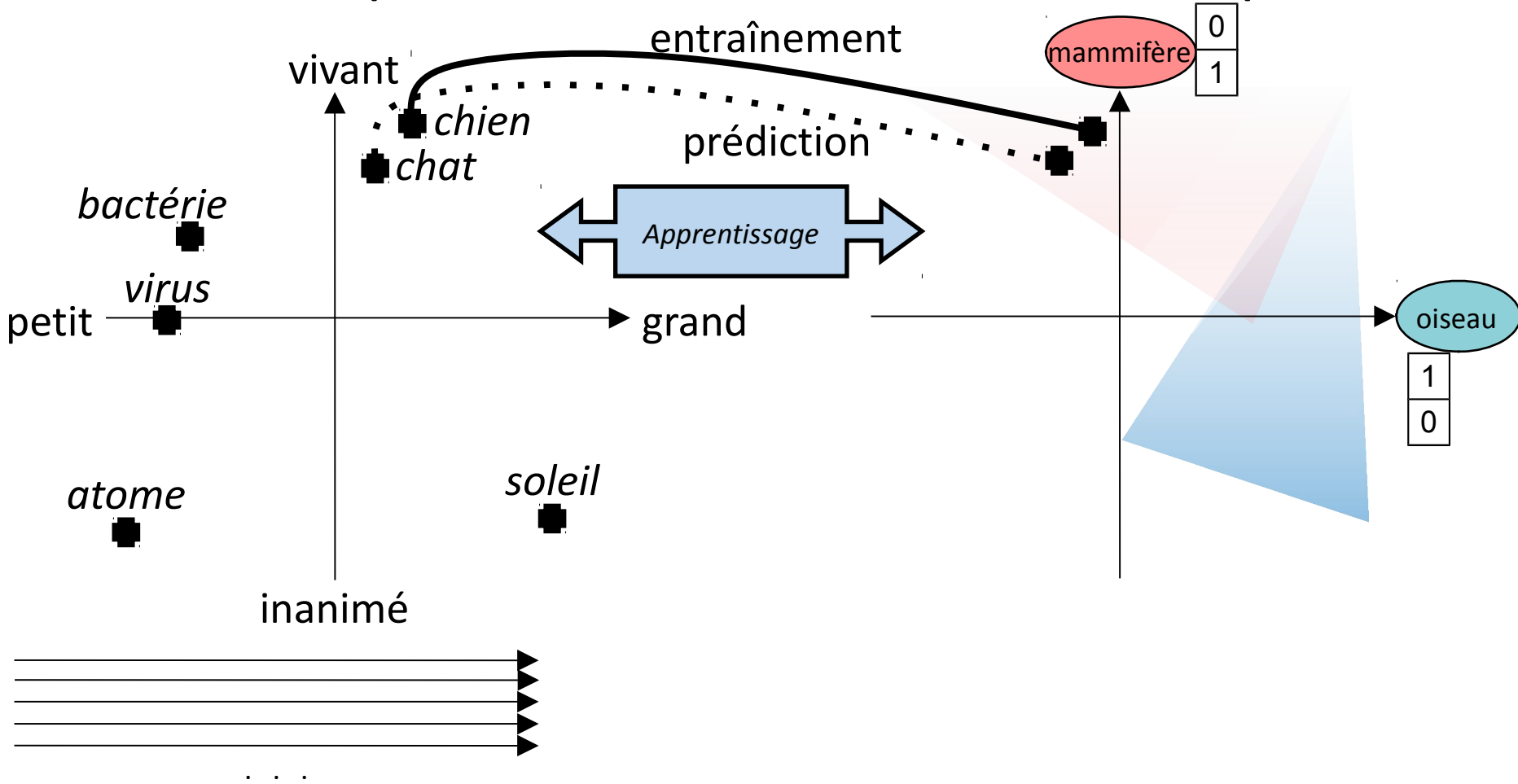
Inconvénients :

- Très grande dimension des vecteurs (performance...)
- Apprentissage « par-cœur »

MÉTHODE : Représentation vectorielle sémantique



MÉTHODE : Représentation vectorielle sémantique



MÉTHODE : Sémantique distributionnelle

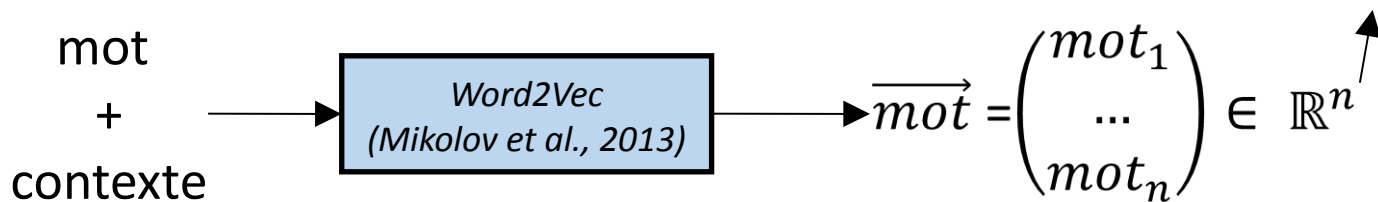
Idée 1 : si on connaît ses contextes d'apparition, on comprend le sens d'un mot

Exemple : « Un chien ronge un _____. »

Idée 2 : si 2 mots partagent les mêmes contextes, ils auront un sens proche

Méthodes permettant de vectoriser le sens d'un mot :

n est un paramètre à choisir (souvent performant entre 100 et 300)

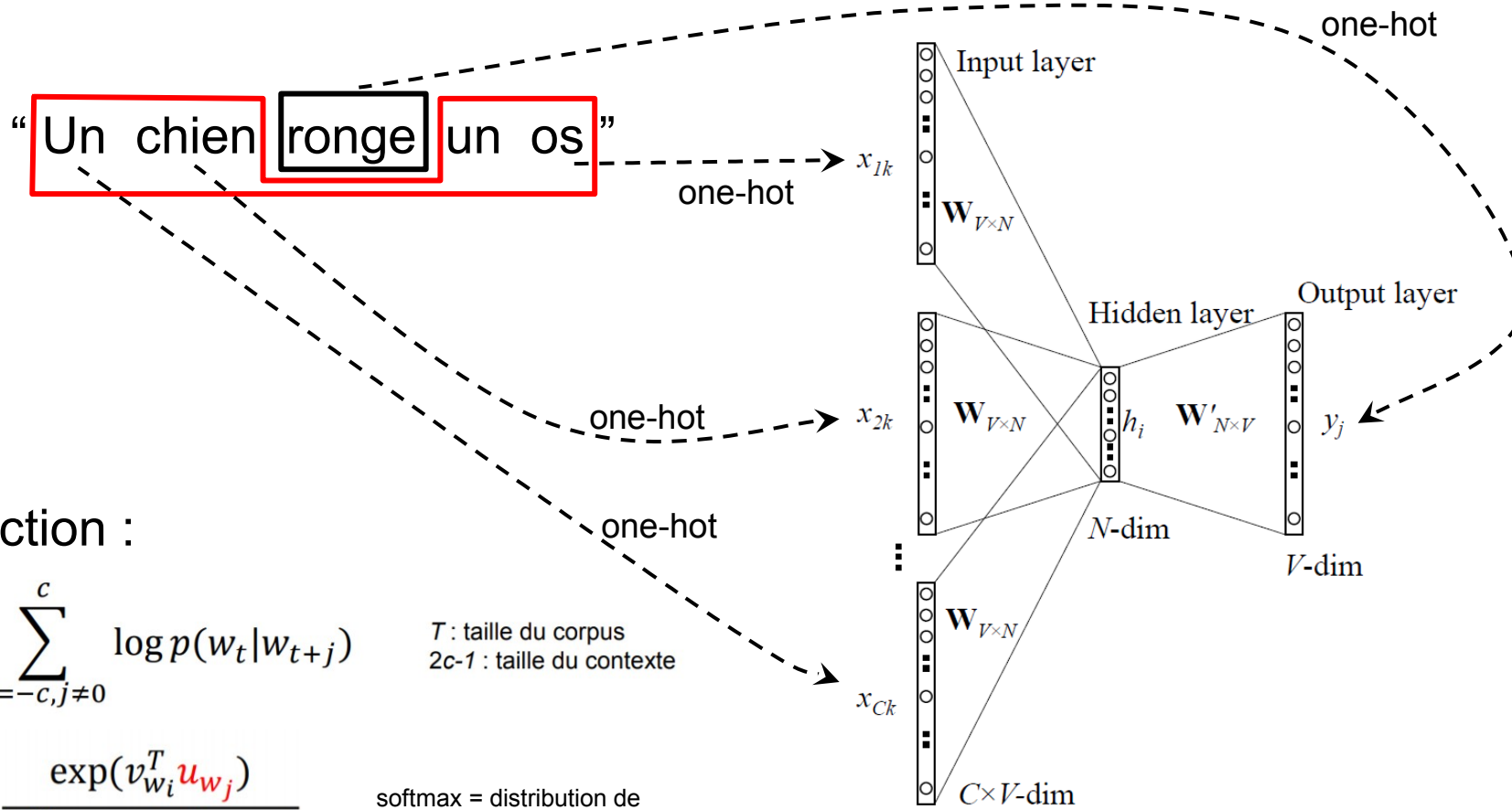


• Proximité spatiale \approx similarité sémantique

• Compositionnalité :

$$\vec{roi} - \vec{homme} + \vec{femme} = \vec{reine} \dots$$

MÉTHODE : Sémantique distributionnelle



Cost function :

$$J_{\theta} = \frac{1}{T} \sum_{t=1}^T \sum_{j=-c, j \neq 0}^c \log p(w_t | w_{t+j})$$

T : taille du corpus
 $2c-1$: taille du contexte

$$p(w_i | w_j) = \frac{\exp(v_{w_i}^T u_{w_j})}{\sum_{k=1}^N \exp(v_{w_k}^T u_{w_j})}$$

softmax = distribution de probabilité sur le vocabulaire

MÉTHODE : Sémantique distributionnelle

cell	Similarité
<i>HCE cell</i>	0.9999
<i>13C-labeled cell</i>	0.9998
<i>parietal cell</i>	0.9989
<i>Schwann cell</i>	0.9965
<i>CD8+ T cell</i>	0.9770
<i>PMN cell</i>	0.9669
<i>macrophage cell</i>	0.9473

seawater	Similarité
<i>sediments</i>	0.7696
<i>sediment sample from a disease free</i>	
<i>fish farm</i>	0.7499
<i>fish farm sediments</i>	0.7342
<i>subterranean brine</i>	0.7320
<i>lagoon on the outskirts of the city of Cagliari</i>	0.7128
<i>petroleum reservoir</i>	0.7095
<i>marine environments</i>	0.7077
<i>marine bivalves</i>	0.6896
<i>sediment samples from diseased farms</i>	0.6870
<i>urine sediments</i>	0.6819
<i>petroleum</i>	0.6576
<i>subterranean environment</i>	0.6497
<i>fresh water</i>	0.6494
<i>fresh water supply</i>	0.6395
<i>Seafood</i>	0.6390
<i>marine</i>	0.6366

MÉTHODE : Sémantique distributionnelle

Avantages :

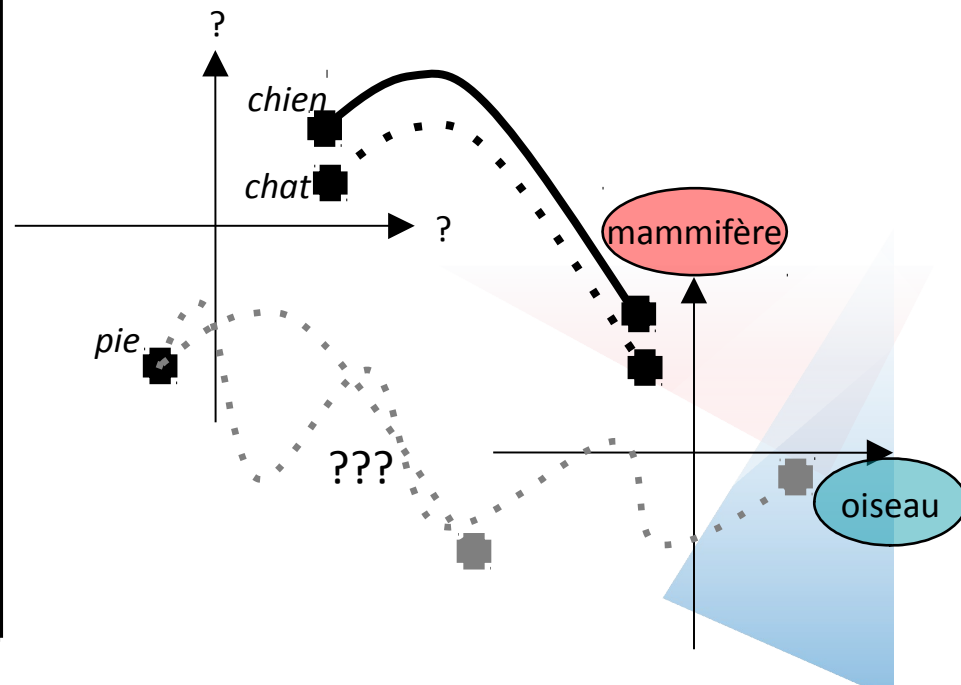
- Représentation vectorielle = nécessaire pour la plupart des algorithmes d'apprentissages
- Permet de faire des prédictions pertinentes de normalisation pour des termes non-rencontrés pendant l'entraînement
- Dimension des vecteurs plus dense (< 1000)

Inconvénient :

- Axes non-interprétables

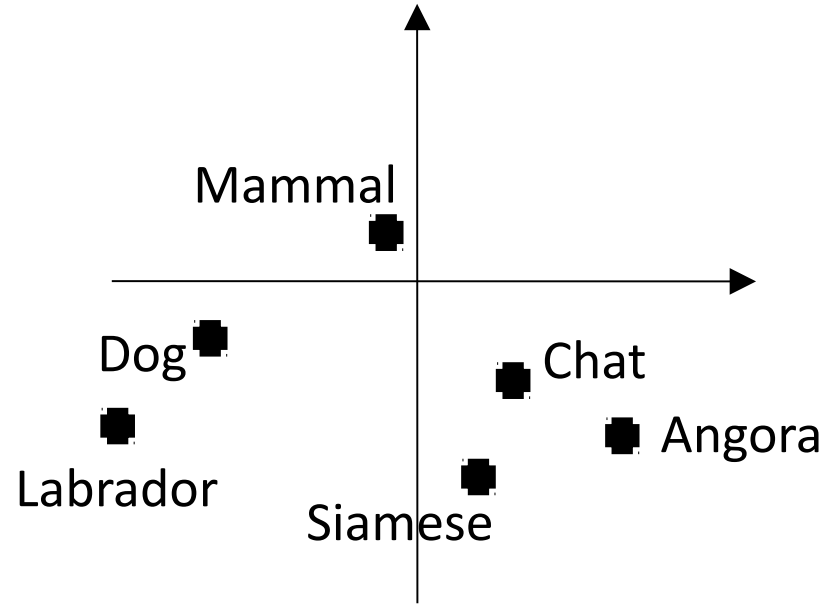
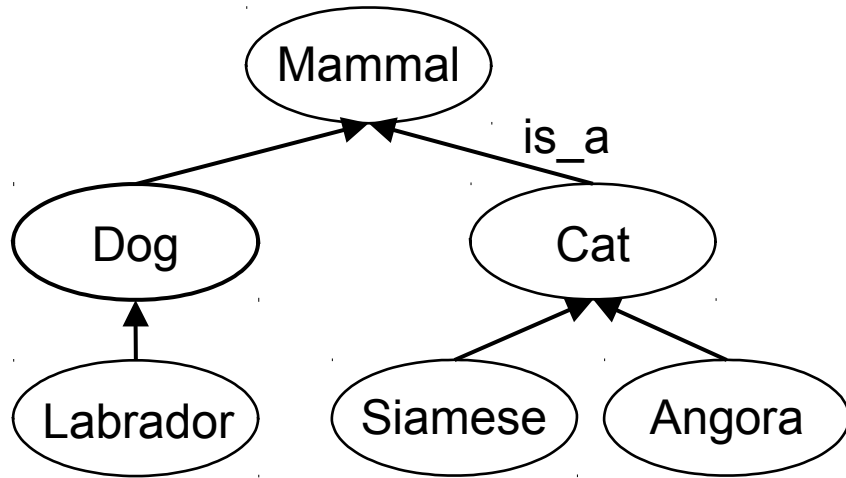
Question :

Comment obtenir des prédictions pertinentes dans le cas où les concepts à prédire n'ont pas été rencontrés lors de l'entraînement ?



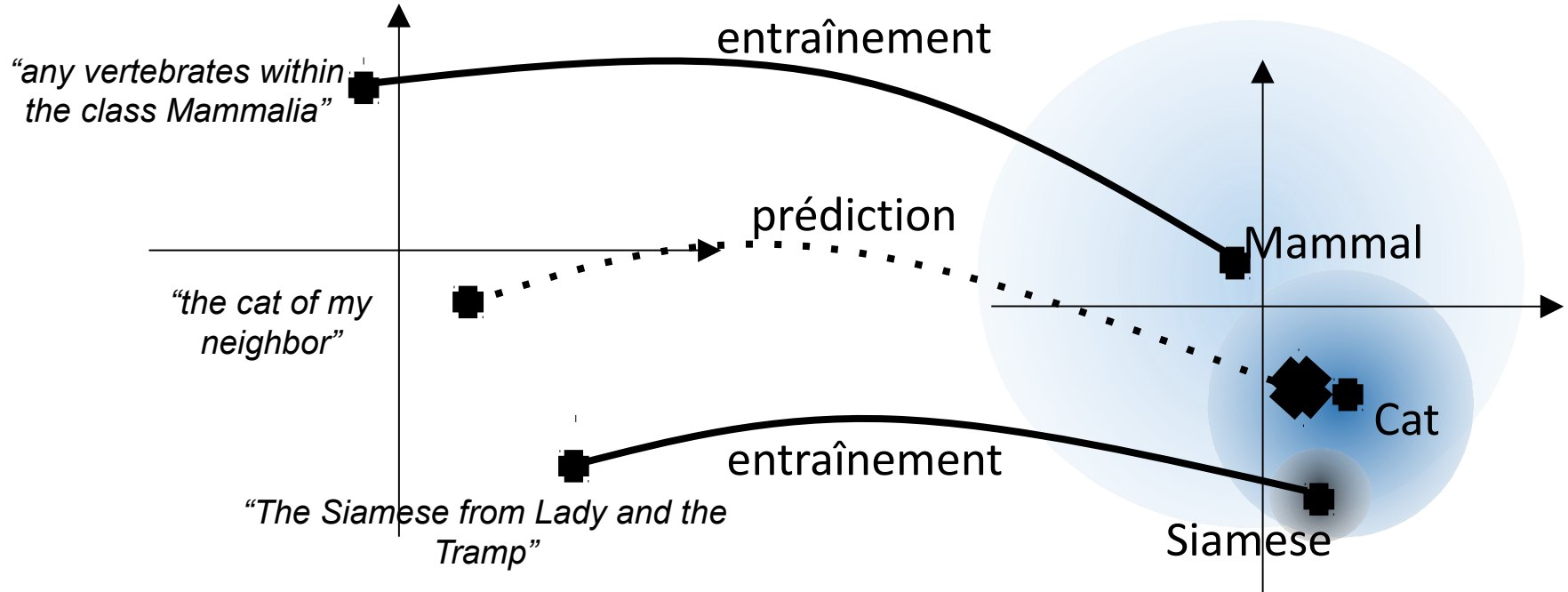
MÉTHODE : Représentation sémantique de concept

Idée : Adopter une représentation vectorielle sémantique pour les concepts également



MÉTHODE : Représentation sémantique de concept

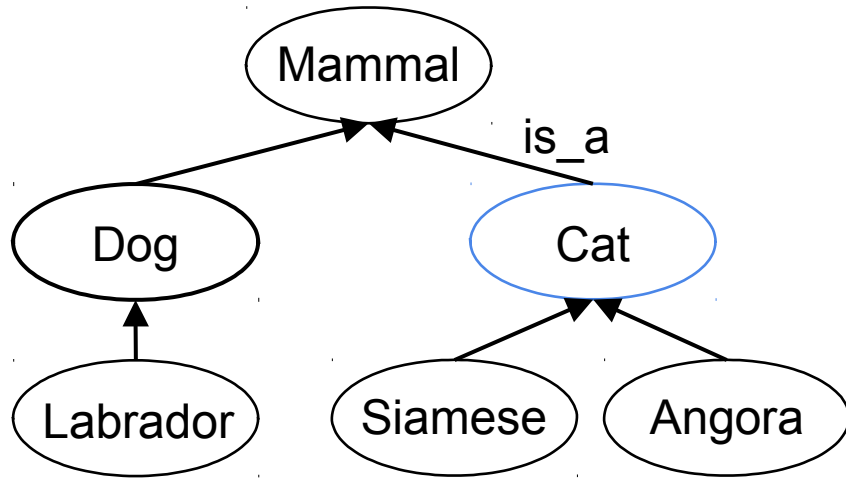
Idée : Adopter une représentation vectorielle sémantique pour les concepts également



Comment peut-on construire de tels vecteurs ?

Quelle information est la plus exploitable ?

MÉTHODE : Représentation sémantique de concept



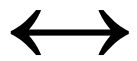
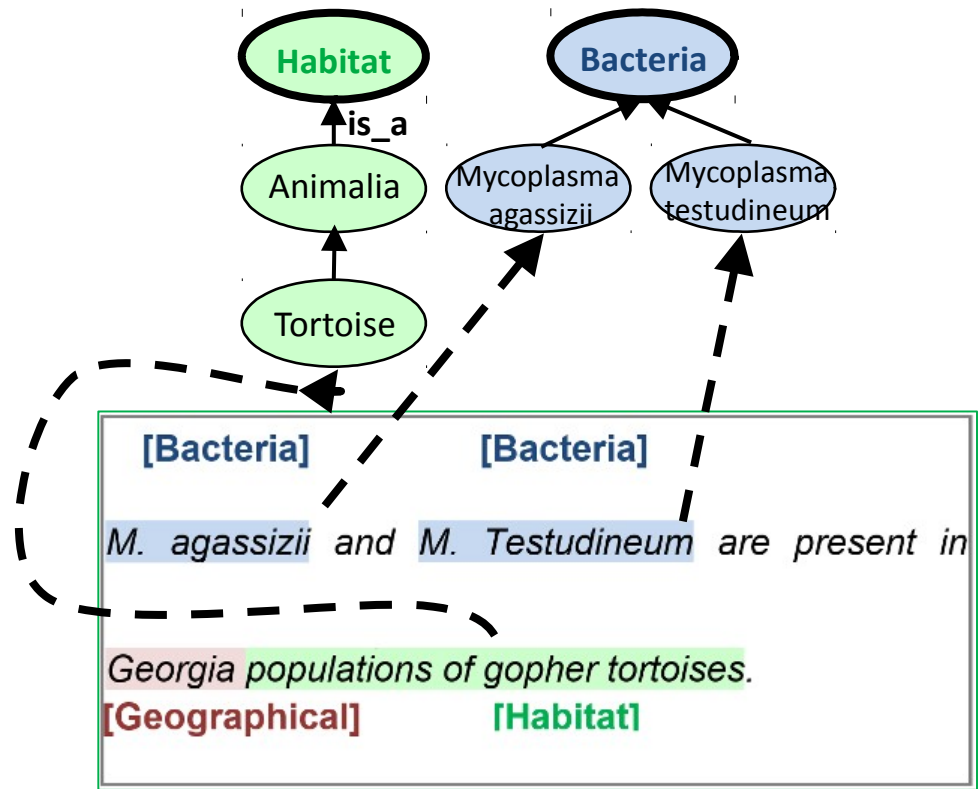
$$\forall k \in \llbracket 1, n \rrbracket, v_{c_k} = (w_{c_k}^0, \dots, w_{c_k}^i, \dots, w_{c_k}^n)$$

$$w_{c_k}^i = \begin{cases} 1 & \text{si } i = k \\ 1 & \text{si } c_i \text{ ancêtre de } c_k \\ 0 & \text{sinon} \end{cases}$$

$$\text{Ex : } v_{\text{Cat}} = [1, 0, 1, 0, 0, 0]$$

Cat	similarité
Cat	1,00
Siamese	0.82
Angora	0.82
Mammal	0.71
Dog	0.50
Labrador	0.41

RÉSULTATS :



Mesure de similarité de Wang et al. (2007)

RÉSULTATS :



ToMap

(Golik Wiktoria, Warnier Pierre, Nédellec Claire, 2011)

Normalise plutôt très bien les termes qui possèdent une tête sémantique de forme similaire à celle d'un label ($\approx 50\%$), mais pas du tout ceux qui n'en possèdent pas.

VS



CONTES

(Ferré Arnaud, Zweigenbaum Pierre, Nédellec Claire, 2017)

Normalise l'ensemble des termes plutôt bien et en particulier ceux ne possédant pas de similarité de forme avec les labels, mais est moins précis pour les termes qui en possèdent.

RÉSULTATS

intégration



Equipe	Score
HONOR	0.73
ToMap (avec règles)	0.66
Turku (2017)	0,63
BOUN (2016)	0,62
ToMap (sans règles)	0.61
CONTES (2017)	0,61
LIMSI (2016)	0,44
Baseline	0,32

● Combining rule-based and embedding-based approaches to normalize textual entities with an ontology
A Ferré, L Deléger, P Zweigenbaum, C Nédellec - LREC 2018

● Representation of complex terms in a vector space structured by an ontology for a normalization task
A Ferré, P Zweigenbaum, C Nédellec - BioNLP 2017

CONCLUSION

Méthode encourageante pour normaliser avec :

- Peu de données d'entraînement
- Un potentiel à répondre au problème de variabilité de forme
- Seulement 2 types d'informations utilisées
(sémantique distributionnelle + information hiérarchique)
- Possibilité de rendre plus interprétable les espaces générés

-> Méthode pouvant servir à peupler automatiquement une ontologie à partir de textes

Normalisation

Méthode :
- par des concepts
d'une ontologie
- applicable à
tout domaine

PERSPECTIVES

- Densifier les vecteurs de concepts (UMLS > 3 millions de concepts)
- Intégrer plus d'information linguistique (notamment syntaxique)
- Quid d'une fonction d'apprentissage non-linéaire ?
- Évoluer vers une méthode non-supervisée qui n'a pas besoin de reconnaissance d'entités nommées

-> S'appuyer sur cette approche pour enrichir (instances mais aussi concepts et relations) automatiquement une ontologie à partir de texte

- Merci de votre attention.
- Des questions?

