

ROCED

Fouille de règles différentielles causales dans les graphes de connaissances

Lucas Simonne ^{*}, Nathalie Pernelle ^{**}, Fatiha Saïs ^{*}

^{*} LRI, Université Paris Saclay

^{**} LIPN, Université Sorbonne Paris Nord



Agenda

Introduction

État de l'Art et Définitions

Première Approche et Résultats

Seconde Approche et Résultats

Règles causales dans les graphes de connaissances

Cadre

- **Causalité** : étude de cause à effet
- Règles causales dans les **graphes de connaissances** - à ne pas confondre avec l'étude des règles d'association dans ce domaine (multiples approches : AMIE, RuDiK, etc.)

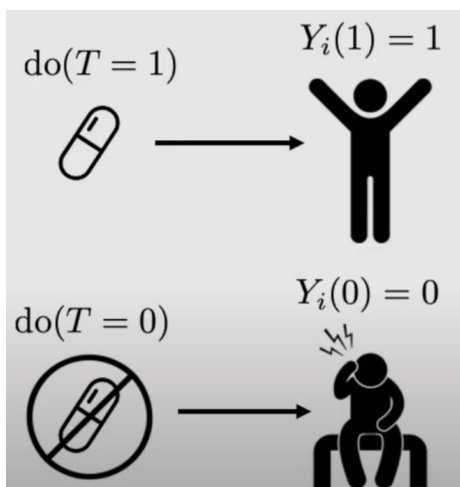
Objectifs

- Proposer un cadre d'étude et une approche originale sur les **graphes de connaissances**
- Expliquer des **différences** sur des propriétés d'une **ontologie** sous forme de **règles** - nous ne sommes pas dans une logique de **prédiction**
- **Intérêt** : découverte de connaissances, aide à la prise de décision

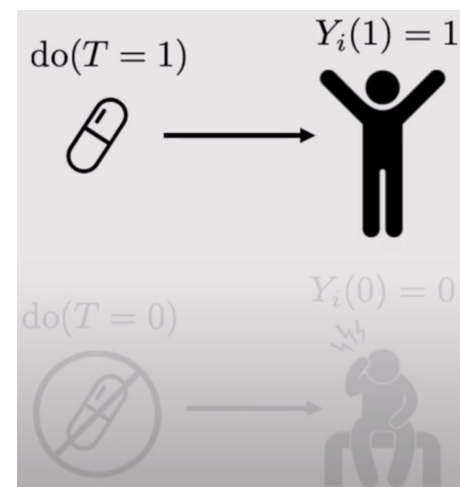
Causalité : Notions (I/II)

Problème Fondamental de la Causalité

- Soit X un **traitement** binaire et Y un **résultat** binaire
 Pour conclure sur l'effet de X sur Y, l'idéal est de comparer pour chaque **unité** les valeurs de Y ayant reçu X et $\neg X$
 Cependant, une unité ne peut recevoir X et $\neg X$, on parle de **contrefactuel**
- **Exemple** : prise d'un traitement et observation de la santé d'un patient



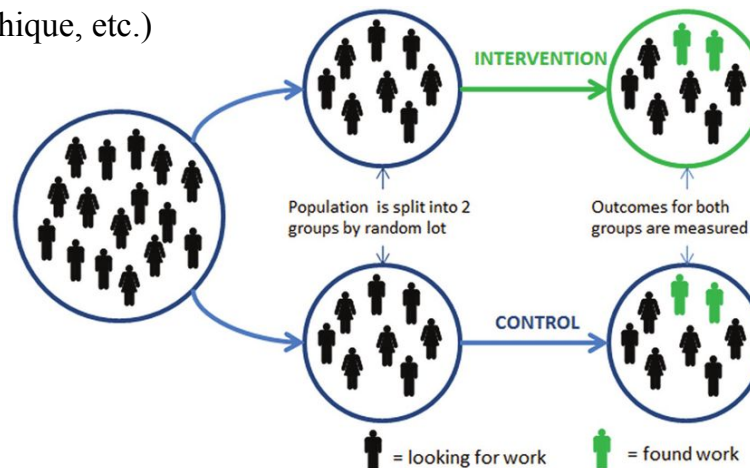
Causal effect
 $Y_i(1) - Y_i(0) = 1$



Causalité : Notions (II/II)

De multiples cadres d'étude existent

- Philosophes (Platon, Aristote) : principe de causalité est un objet d'étude important
- Fin XXe : Rubin et Holland - **modèles à résultats potentiels**
- Début XXIe : Pearl - **structures de graphes causaux**
- Autres approches : réseaux bayésiens, réseaux profonds, etc.
- Standard : **attribution aléatoire** des traitements (également appelé test A/B), mais utilisation **limitée** (coût, éthique, etc.)



Haynes et al. (2012)

Nous nous situons dans le cadre d'étude des résultats potentiels dans les données observationnelles.

Peu d'approches combinent causalité et graphe de connaissances

Approches récentes sur des Graphes de Connaissances

- Munch et al. (2019) transforme les données RDF en **schéma relationnel** puis entraîne un **réseau bayésien** pour obtenir des distributions conjointes.
- Limites :
 - Variables catégorielles seulement - pas de gradualité

Première approche à étudier la causalité dans les graphes de connaissances dans le cadre des résultats potentiels.

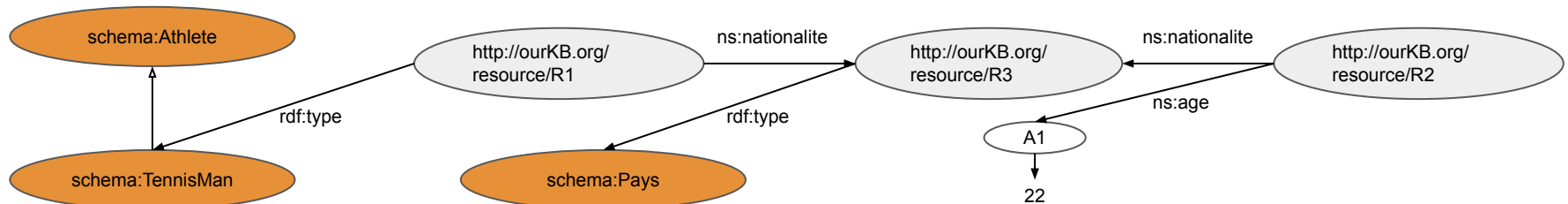
Approches récentes en Causalité dans le cadre des résultats potentiels

- Stuart (2010) : vue d'ensemble des méthodes utilisées pour les résultats potentiels
 - **appariement** exact, tronqué, pondéré, stratifié.
- Beaucoup d'approches cherchent à réduire la complexité du problème :
 - Li et al. (2013) : règles d'associations causales à partir de règles d'associations
 - Li et al. (2017) : arbre causal - chemins causaux sur critère modifié

Définitions Clés (I/III)

Graphe de Connaissances

- Ensemble de **prédicats** < sujet, propriété, objet >
 - sujets et certains objets : URIs correspondant aux noeuds
 - propriétés : URIs correspondant aux arêtes
 - autres objets : littéraux (texte, image, entier, etc.)
- **2 niveaux d'abstractions** : les données (RDF) et ontologie (OWL, RDFS)



Définitions Clés (II/III)

Règle Différentielle Causale

- Soit $C \in C$ la classe cible de l'ontologie et pr le résultat étudié, une **règle différentielle causale** $B \Rightarrow R(X, Y)$ est définie par :

$$\begin{aligned}
 & C(X_1) \wedge C(X_2) [\wedge ST_i(X_1) \wedge ST_i(X_2)] \\
 & [\bigwedge_{k=1; k=K} pt_k(X_1, V_1) \wedge pt_k(X_2, V_2) \wedge compare_t(V_1, V_2)] \\
 & [\bigwedge_{l=1; l=L} pt_l(X_1, v_1) \wedge pt_l(X_2, v_2)] \wedge po(X_1, U_1) \wedge po(X_2, U_2) \\
 & \Rightarrow compare_o(U_1, U_2)
 \end{aligned}$$

où :

- ST_i une **strate** optionnelle
- pt_k (pt_l) : **chemin de propriétés** menant vers un **traitement** (valeur numérique ou ordonnées)
- $compare$: opérateurs arithmétiques $<, >$
- po : **chemin de propriétés** menant vers un **résultat** (numérique ou ordonné)
- ex : $TennisMan(i_1) \wedge TennisMan(i_2) \wedge age(i_1, v_1) \wedge age(i_2, v_2) \wedge inferior(v_1, v_2) \wedge performance(i_1, z_1) \wedge performance(z_2, z_2) \Rightarrow superior(z_1, z_2)$

Une telle règle indique une relation entre une différence de traitements et une différence de résultats, applicable pour des éléments d'une strate.

Définitions Clés (III/III)

Strate

- Soit i_1 une instance de la classe cible C .
Une **strate** $ST_j(i_1)$ est une conjonction de prédicats qui correspond à un motif de graphe en RDF, dont la racine est i_1 , et pour lequel les feuilles représentent des valeurs littérales ou des classes.
Un tel motif peut être limité par une profondeur d correspondant à la longueur maximale d'un chemin atteint par ce motif de graphe.
- exemples :
 - $TennisMan(X)$
 - $TennisMan(X) \wedge nationalite(X, "France")$

Quels critères pour les règles ?

Métrique

$$\vec{B} \wedge op(V_1, V_2) \Rightarrow op(U_1, U_2)$$

$$causal_r(R) = \frac{supp(\vec{B} \wedge op(V_1, V_2) \Rightarrow op(U_1, U_2))}{supp(\vec{B} \wedge op(V_1, V_2) \Rightarrow op(U_2, U_1))}$$

- Inspirée de Li et al. (2012)
- Soit R une règle maximale spécifique, R est **cohérente** si $causal_r(R)$ est significativement supérieure à 1.

Règle Cohérente

Une règle $R_1 : \vec{ST}_1 \wedge \vec{T} \Rightarrow compare(U_1, U_2)$ est cohérente ssi : $\forall ST_2 \sqsubseteq ST_1, R_2 : \vec{ST}_2 \wedge \vec{T} \Rightarrow compare(U_1, U_2)$ est aussi cohérente.

Une strate ST_1 est maximale spécifique ssi il n'existe pas $ST_2 \sqsubseteq ST_1$ telle que $depth(ST_2) \leq d$.

Pourquoi étudier les règles les plus spécifiques ?

Exemple de paradoxe

		Condition		
		Mild	Severe	Total
Traitement	A	15% (210/1400)	30% (30/100)	16% (240/1500)
	B	10% (5/50)	20% (100/500)	19% (105/550)

- Inversion des effets : quel traitement choisir ?

Conséquences

- Pour éviter inversement des effets : on garde les règles spécifiques (règle générale est ok si ses spécifiques sont ok).
Stratégie avec contraintes (impose que ces strates soient instanciées).

Algorithme

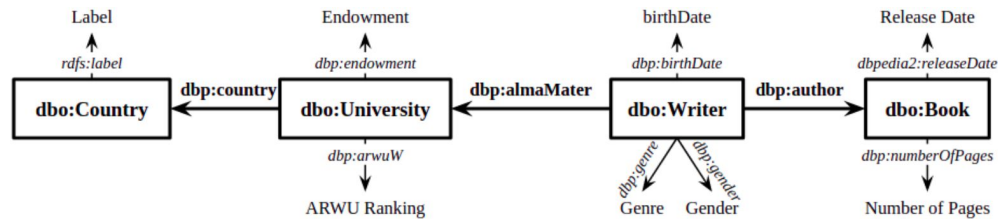
- Entrées : Graphe de connaissance KG , Classe cible C ; seuil de support $minsup$; paramètre de test statistique x ; ensemble vide des règles causales CR ; nombre de traitement n_t
- (1) Construction des strates maximales spécifiques dont le support $> minsup$ par combinaison de strates élémentaires
- (2) Construction des règles par sélection de strates identiques aux traitements près (n_t)
- (3) Calcul de $causal_{\tau}$ pour tous les R : si la borne minimale de $seb-kg(R)$ de l'intervalle de confiance à $x\% > 1$, ajout à CR .
- (4) Généralisation de règles par fusion - Algorithme récursif

L'algorithme a été testé sur 2 jeux de données.

Données DBPedia - Expérimentations

Jeu de données utilisé

- Jeu de données utilisé par Munch et al. (2019) : seule approche combinant causalité et graphe de connaissances



6908 triplets
 Classe Cible : Auteur
 Intervalle de confiance à
 90%

Résultats

- Confirment les résultats obtenus dans Munch et al. (2019)
- 11 règles spécifiques (**graduelles**) avec comme traitements et résultats :
 - **Plus** l'université de l'auteur est mieux classée, le **plus** tôt il publiera son premier livre.
 - **Plus** l'auteur est né tôt, le **plus tard** il publiera son premier livre
- Ces traitements sont appliqués à différentes strates :
 - auteurs nés aux USA ;
 - auteurs nés à une certaine époque ;
 - etc.

Limites et Extension

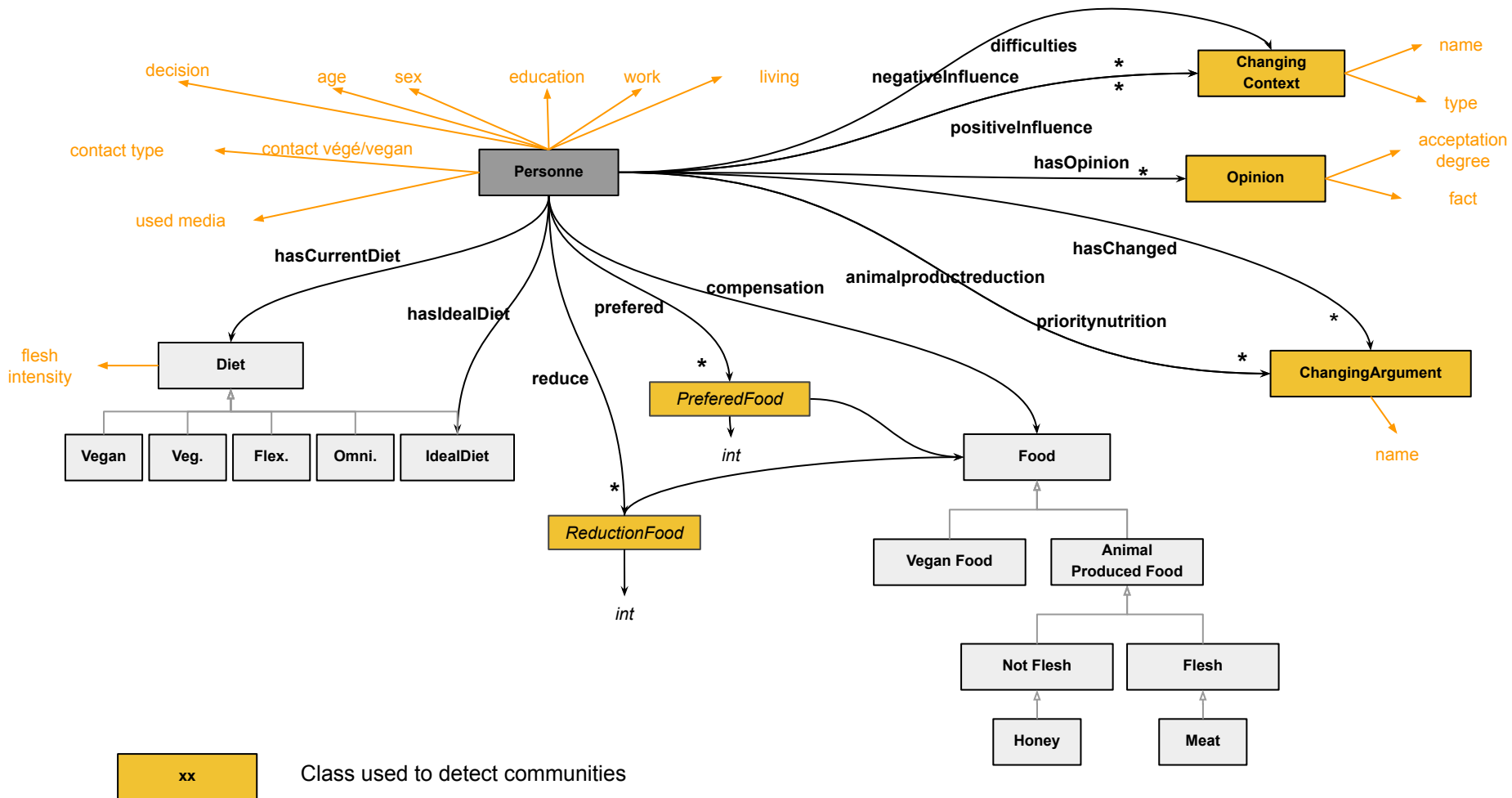
Limites

- Nombre de strates augmente avec la complexité de l'ontologie : instances dispersées, compliqué d'obtenir des métriques statistiquement significatives
- Des propriétés peuvent être non fonctionnelles (par exemple : un auteur peut avoir écrit plus d'un livre)

Extension

- Supprimer des propriétés :
 - Utiliser la sémantique de l'ontologie pour ne pas considérer des strates qui n'ont pas de sens
 - L'expert peut choisir de **ne pas étudier** certaines propriétés non pertinentes
- Simplifier la représentation des strates :
 - certaines propriétés peuvent être **étudiées conjointement**
 - l'expert définit les ensembles à étudier ensemble
 - une mesure de similarité est définie pour chaque paire de propriétés de cet ensemble (**co-occurrence**)
 - détection de **communautés** pour générer des prédicats abstraits

Ontologie du 2e jeu de données : Vitamin



Données Vitamin - Expérimentations

Jeu de données utilisé

- Vitamin : > 91k prédicats
- Résultat :
 - Classe cible : Personne
 - Différence entre le **régime actuel** et le **régime idéal**

EXTRACT OF COMMUNITIES - *Vitamin*

Community 1
Livestock is a major contributor to global warming
Vegetarian diets are healthier
Animals suffer
The problem is not breeding but industrial breeding
Community 2
Vegan diets are deficient
Human Nature is to eat animals
Vegetarian diets are deficient
Eating animal products makes me happy

Résultats

- Règles obtenues :
 - Règle la plus générale : être sensible au bien être animal et au bien être animal par rapport à ne pas considérer ces sujets peut expliquer une différence de changement de régime.
 - Règles plus spécifiques également obtenues : être une femme par rapport à être homme (pour les individus vivant en ville)
- Evaluation des experts :
 - Détection de communautés : intéressante pour exprimer des communautés polarisées ayant une sémantique
 - Règles : règles obtenues intuitives et faciles à interpréter quand les strates sont générales

Conclusion et Perspectives

Conclusion

- Une première approche de détection de règles différentielles causales adaptée au KG
- Premières expérimentations montrent que :
 - ontologies simples : appariement et strates exploitent l'ensemble du KG
 - ontologies plus complexes : stratégie de génération de propriétés abstraites basée sur détection de communauté permet de prendre en compte les ontologie plus complexes

Travaux futurs

- Utilisation de l'ontologie pour guider la construction des règles
- Nouvelle approche reposant sur des plongements de graphes pour déterminer des effets moyens

Merci pour votre attention.

simonne@lri.fr