

# Assister l'édition manuelle de données RDF à l'aide du raisonnement à partir de cas

*Journées IC 2021*

Nicolas Lasolle, Olivier Bruneau, Jean Lieber,  
Emmanuel Nauer, Siyana Pavlova

2 juillet 2021



Archives Henri Poincaré  
Philosophie et Recherches sur les  
Sciences et les Technologies



**IMPACT**  
**OLKi**

## Web sémantique

- ▶ Données structurées selon le modèle RDF
- ▶ Connaissances du domaine représentées en RDFS
- ▶ Graphes RDF interrogeables et modifiables grâce à SPARQL



## Application pour un corpus historique

Les travaux présentés sont appliqués pour enrichir et exploiter le corpus de la correspondance d'Henri Poincaré.



# Problématique

## Constat

L'édition manuelle de données RDF est une tâche fastidieuse avec un risque d'erreurs important :

- ▶ Erreur de duplication ;
- ▶ Erreur d'ambiguïté ;
- ▶ Erreur de frappe.

# Problématique

## Constat

L'édition manuelle de données RDF est une tâche fastidieuse avec un risque d'erreurs important :

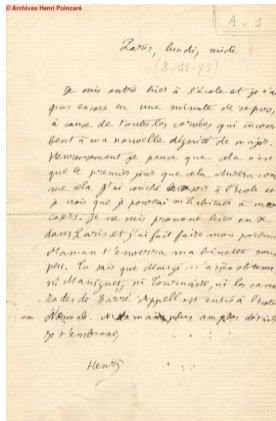
- ▶ Erreur de duplication ;
- ▶ Erreur d'ambiguïté ;
- ▶ Erreur de frappe.

## Objectif

Proposer, implémenter, et évaluer des méthodes pour assister les contributeurs dans cette tâche.

# Le corpus de la correspondance d'Henri Poincaré

# Le corpus de la correspondance d'Henri Poincaré

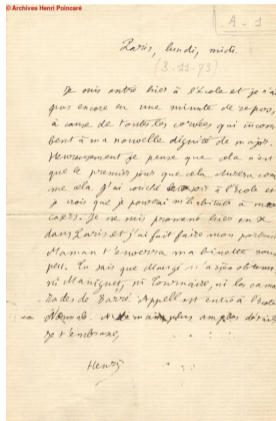


Henri Poincaré à Aline  
Poincaré, 1873

## Un corpus historique

- ▶ Plus de 2000 lettres (échanges d'ordre scientifique, administratif et privé)

# Le corpus de la correspondance d'Henri Poincaré

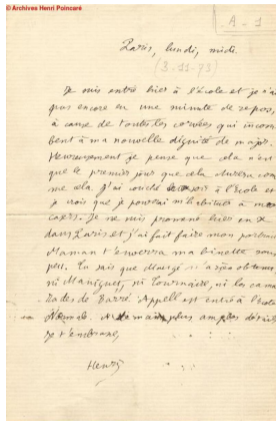


Henri Poincaré à Aline  
Poincaré, 1873

## Un corpus historique

- ▶ Plus de 2000 lettres (échanges d'ordre scientifique, administratif et privé)
- ▶ Source d'informations importante pour les historiens

# Le corpus de la correspondance d'Henri Poincaré



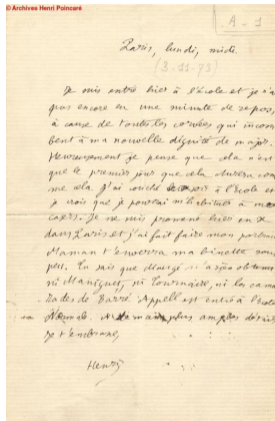
Henri Poincaré à Aline  
Poincaré, 1873

## Un corpus historique

- ▶ Plus de 2000 lettres (échanges d'ordre scientifique, administratif et privé)
- ▶ Source d'informations importante pour les historiens
  - Théories scientifiques



# Le corpus de la correspondance d'Henri Poincaré



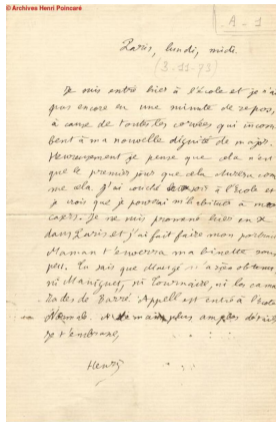
Henri Poincaré à Aline  
Poincaré, 1873

## Un corpus historique

- ▶ Plus de 2000 lettres (échanges d'ordre scientifique, administratif et privé)
- ▶ Source d'informations importante pour les historiens
  - Théories scientifiques
  - Institutions et sociétés savantes



# Le corpus de la correspondance d'Henri Poincaré



Henri Poincaré à Aline  
Poincaré, 1873

## Un corpus historique

- ▶ Plus de 2000 lettres (échanges d'ordre scientifique, administratif et privé)
- ▶ Source d'informations importante pour les historiens
  - Théories scientifiques
  - Institutions et sociétés savantes
  - Contexte politique, social et culturel
- ▶ L'édition de ce corpus est un projet collectif mené par les Archives Henri-Poincaré

# Le corpus de la correspondance d'Henri Poincaré

L'édition numérique et le site Omeka S



## LA CORRESPONDANCE D'HENRI POINCARÉ

[Accueil](#)

[Lettres](#)

[Recherche](#)

[Index](#)

### LETTRE : Henri Poincaré à Aline Boutroux - 3 novembre 1873

[Modifier l'item](#)

[Transcription](#)

[Métadonnées](#)

[Citer ce document](#)

[3 novembre 1873<sup>1</sup>]

Paris, lundi, midi.

Je suis entré hier à l'école<sup>2</sup> et je n'ai pas encore eu une minute de repos, à cause de toutes les nouvelles corvées qui incombent à ma nouvelle dignité de **major**. Heureusement je pense que cela n'est que le premier jour que cela durera comme cela. J'ai couché hier soir à l'École et je crois que je pourrai m'habituer à mon **casern**. Je me suis promené hier en **X**<sup>3</sup> dans Paris et j'ai fait faire mon portrait<sup>4</sup>. **Maman** t'enverra ma binette sous peu. Tu sais que **Maugé**<sup>5</sup> n'a rien obtenu, ni **Maniguet**<sup>6</sup>, ni **Tournaire**, ni les camarades de **Barré**<sup>7</sup>. **Appell** est entré à l'École Normale<sup>8</sup>. À demain plus amples détails.

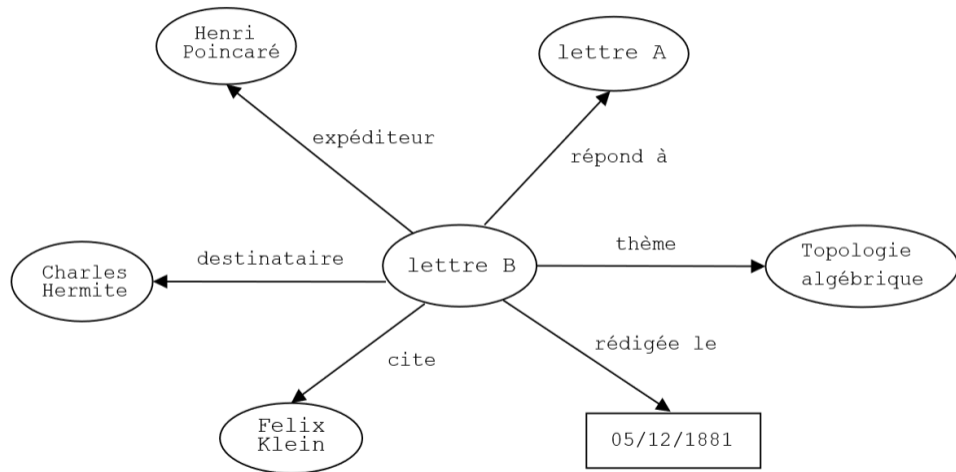
Je t'embrasse,

**Henri**

1. Comme de nombreuses autres, cette lettre n'avait pas de destinataire précis. La référence à **Eugénie Poincaré** laisse supposer qu'elle s'adressait plutôt à la sœur de **Poincaré**, **Aline**. On sait que toute la famille avait accompagné le jeune polytechnicien à Paris pour sa rentrée. Il est probable qu'**Eugénie Poincaré** était restée quelques jours à Paris tandis qu'**Aline** et **Émile-Léon Poincaré** étaient retournés à Nancy.↩

# Le corpus de la correspondance d'Henri Poincaré

## Resource Description Framework (RDF)



# Le corpus de la correspondance d'Henri Poincaré

## SPARQL

Ce langage permet de formuler des requêtes afin d'explorer le graphe RDF du corpus

### Exemple de requête (informelle)

$Q =$  | Donner les lettres envoyées par Henri Poincaré  
à Félix Klein entre 1890 et 1900 et ayant  
pour thème la géométrie

# Des méthodes pour assister l'édition manuelle de données RDF

# Comment assister l'édition manuelle de données RDF ?

## Problématique

L'édition manuelle de données (i.e. saisir les triplets) est souvent une tâche longue et fastidieuse.



# Comment assister l'édition manuelle de données RDF ?

## Problématique

L'édition manuelle de données (i.e. saisir les triplets) est souvent une tâche longue et fastidieuse.

## Objectif

Fournir une liste de suggestions ordonnées pour assister les contributeurs dans cette tâche.

# Comment assister l'édition manuelle de données RDF ?

## Problématique

L'édition manuelle de données (i.e. saisir les triplets) est souvent une tâche longue et fastidieuse.

## Objectif

Fournir une liste de suggestions ordonnées pour assister les contributeurs dans cette tâche.

## Questions d'édition

$\langle \boxed{?s} \ p \ o \rangle$        $\langle s \ \boxed{?p} \ o \rangle$        $\langle s \ ?p \ \boxed{?o} \rangle$        $\langle s \ \boxed{?p} \ ?o \rangle$       ...

# Comment assister l'édition manuelle de données RDF ?

## Problématique

L'édition manuelle de données (i.e. saisir les triplets) est souvent une tâche longue et fastidieuse.

## Objectif

Fournir une liste de suggestions ordonnées pour assister les contributeurs dans cette tâche.

## Questions d'édition

$\langle \boxed{?s} \ p \ o \rangle$      $\langle s \ \boxed{?p} \ o \rangle$      $\langle s \ ?p \ \boxed{?o} \rangle$      $\langle s \ \boxed{?p} \ ?o \rangle$     ...

## Exemple

$\langle \text{lettreA} \ \text{destinataire} \ \boxed{?o} \rangle$

## Proposition d'un outil d'édition

4 versions du système de suggestions :

**basique** classe les suggestions selon l'ordre alphabétique ;

**déductif** bénéficie de l'utilisation des connaissances de l'ontologie ;

**à base de cas** s'appuie sur des problèmes d'édition similaires au problème courant ;

**combiné** combine les deux précédentes méthodes.

## Proposition d'un outil d'édition

4 versions du système de suggestions :

**basique** classe les suggestions selon l'ordre alphabétique ;

**déductif** bénéficie de l'utilisation des connaissances de l'ontologie ;

à **base de cas** s'appuie sur des problèmes d'édition similaires au problème courant ;

**combiné** combine les deux précédentes méthodes.

# Proposition d'un outil d'édition

4 versions du système de suggestions :

**basique** classe les suggestions selon l'ordre alphabétique ;

**déductif** bénéficie de l'utilisation des connaissances de l'ontologie ;

**à base de cas** s'appuie sur des problèmes d'édition similaires au problème courant ;

**combiné** combine les deux précédentes méthodes.

# Proposition d'un outil d'édition

4 versions du système de suggestions :

**basique** classe les suggestions selon l'ordre alphabétique ;

**déductif** bénéficie de l'utilisation des connaissances de l'ontologie ;

**à base de cas** s'appuie sur des problèmes d'édition similaires au problème courant ;

**combiné** combine les deux précédentes méthodes.

# Assister l'édition manuelle de données RDF

Systeme déductif

## Méthodologie

Le système s'appuie sur les domaines (`rdfs:domain`) et co-domaines (`rdfs:range`) des propriétés de l'ontologie.

Question d'édition `<lettreA destinataire [?]>`

## Solution pour le classement des valeurs potentielles

Comme les classes `Personne` et `Institution` font partie du co-domaine de la propriété `destinataire`, cette connaissance est utilisée pour favoriser les instances de ces classes.



# Assister l'édition manuelle de données RDF

Système déductif

## Méthodologie

Le système s'appuie sur les domaines (`rdfs:domain`) et co-domaines (`rdfs:range`) des propriétés de l'ontologie.

Question d'édition `<lettreA destinataire [?]o>`

## Solution pour le classement des valeurs potentielles

Comme les classes `Personne` et `Institution` font partie du co-domaine de la propriété `destinataire`, cette connaissance est utilisée pour favoriser les instances de ces classes.

## Limite

Plus de 1000 personnes dans le corpus !

# Assister l'édition manuelle de données RDF

Système à base de cas

## Raisonnement à partir du cas

Une méthodologie qui s'appuie sur des situations précédentes pour résoudre des problèmes nouveaux.

# Assister l'édition manuelle de données RDF

Système à base de cas

## Raisonnement à partir du cas

Une méthodologie qui s'appuie sur des situations précédentes pour résoudre des problèmes nouveaux.

## Définitions

Cas :  $(x, y)$  où  $x$  : problème et  $y$  : solution

Problème cible : le problème à résoudre  $x^{\text{cible}}$

Cas source :  $(x^s, y^s) \in \text{BaseCas}$  utilisé pour résoudre  $x^{\text{cible}}$

# Assister l'édition manuelle de données RDF

Système à base de cas

## Le problème cible

$x^{\text{cible}}$  =

<b>question :</b> <lettreA destinataire <input type="text" value="?o"/> >
<b>contexte :</b> <lettreA expéditeur henriPoincaré>
<lettreA thème topologieAlgébrique>
<lettreA cite charlesHermite>
<lettreA rédigéeEn 1880>

# Assister l'édition manuelle de données RDF

Système à base de cas

Le problème cible

$x^{\text{cible}}$  =

<b>question :</b> <code>&lt;lettreA destinataire ?o &gt;</code>
<b>contexte :</b> <code>&lt;lettreA expéditeur henriPoincaré &gt;</code>
<code>&lt;lettreA thème topologieAlgébrique &gt;</code>
<code>&lt;lettreA cite charlesHermite &gt;</code>
<code>&lt;lettreA rédigéeEn 1880 &gt;</code>

$y^{\text{cible}}$  = valeur pour `?o`

# Assister l'édition manuelle de données RDF

Système à base de cas

Le problème cible

$x^{\text{cible}}$  =

<b>question :</b> $\langle$ lettreA destinataire <input type="text" value="?o"/> $\rangle$
<b>contexte :</b> $\langle$ lettreA expéditeur henriPoincaré $\rangle$
$\langle$ lettreA thème topologieAlgébrique $\rangle$
$\langle$ lettreA cite charlesHermite $\rangle$
$\langle$ lettreA rédigéeEn 1880 $\rangle$

$y^{\text{cible}}$  = valeur pour

Base de cas = la base RDF  $\mathcal{D}$

# Assister l'édition manuelle de données RDF

Système à base de cas

Le problème cible

$$x^{\text{cible}} = \begin{array}{l} \text{question : } \langle \text{lettreA destinataire } \boxed{?o} \rangle \\ \text{contexte : } \langle \text{lettreA expéditeur henriPoincaré} \rangle \\ \langle \text{lettreA thème topologieAlgébrique} \rangle \\ \langle \text{lettreA cite charlesHermite} \rangle \\ \langle \text{lettreA rédigéeEn 1880} \rangle \end{array}$$

$$y^{\text{cible}} = \text{valeur pour } \boxed{?o}$$

Base de cas = la base RDF  $\mathcal{D}$

$x^s$  = une ressource de  $\mathcal{D}$

$y^s$  = une valeur candidate pour  $\boxed{?o}$   
 $\langle x^s \text{ destinataire } y^s \rangle$

# Assister l'édition manuelle de données RDF

Système à base de cas

$x^{\text{cible}} =$

**question :** <lettreA destinataire ?o>

<lettreA expéditeur henriPoincaré>

**contexte :** <lettreA thème topologieAlgébrique>

<lettreA cite charlesHermite>

<lettreA rédigéeEn 1880>

## Requête SPARQL initiale

Retrouver les ressources pour lesquelles le contexte est le même que lettreA.



# Assister l'édition manuelle de données RDF

Systeme à base de cas

## Requête initiale

$Q =$  | Donner les lettres envoyées par Poincaré en 1880  
ayant pour thème la topologie algébrique  
et citant Charles Hermite.

# Assister l'édition manuelle de données RDF

Système à base de cas

## Requête initiale

$$Q = \left| \begin{array}{l} \text{Donner les lettres envoyées par Poincaré en 1880} \\ \text{ayant pour thème la topologie algébrique} \\ \text{et citant Charles Hermite.} \end{array} \right.$$

## Problème

Ce n'est pas courant d'avoir deux ressources avec le même contexte. Comment retrouver les cas les plus similaires dans la base ?

# Assister l'édition manuelle de données RDF

Système à base de cas

## Requête initiale

$Q =$  | Donner les lettres envoyées par Poincaré en 1880  
ayant pour thème la topologie algébrique  
et citant Charles Hermite.

## Problème

Ce n'est pas courant d'avoir deux ressources avec le même contexte. Comment retrouver les cas les plus similaires dans la base ?

## Proposition

Utiliser d'un mécanisme de recherche flexible s'appuyant sur des règles de transformation de requêtes SPARQL.

# Assister l'édition manuelle de données RDF

Système à base de cas

## Requête initiale

$Q =$  | Donne moi les lettres envoyées par Poincaré en 1880  
ayant pour thème la **topologie algébrique**  
et citant Charles Hermite.

## Requête générée (en généralisant le thème)

$Q_1 =$  | Donner les lettres envoyées par Poincaré en 1880  
ayant pour thème les **mathématiques**  
et citant Charles Hermite.

# Assister l'édition manuelle de données RDF

Système à base de cas

## Requête initiale

$Q =$  | Donner les lettres envoyées par Poincaré en 1880  
ayant pour thème la topologie algébrique  
et citant **Charles Hermite**.

Requête générée (en remplaçant une personne citée par une personne liée)

$Q_2 =$  | Donner les lettres envoyées par Poincaré en 1880  
ayant pour thème la topologie algébrique  
et citant **Paul Appell**.

# Assister l'édition manuelle de données RDF

Système à base de cas

## Requête initiale

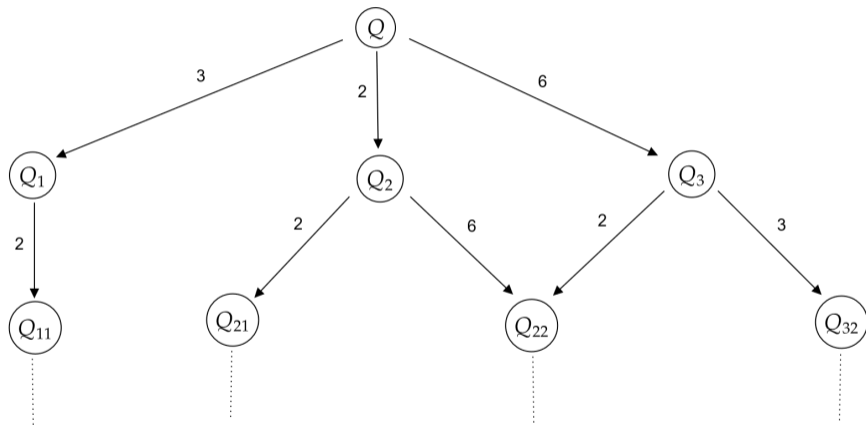
$Q =$  | Donner les lettres envoyées par Poincaré en 1880  
ayant pour thème la topologie algébrique  
et citant Charles Hermite.

Requête générée (en étendant les bornes temporelles liées à la date de rédaction)

$Q_3 =$  | Donner les lettres envoyées par Poincaré entre 1875  
et 1885 ayant pour thème la topologie algébrique  
et citant Charles Hermite.

# Assister l'édition manuelle de données RDF

Coûts de transformation



# Implémentation



# Évaluation des méthodes

# Évaluation humaine

## Méthodologie

- ▶ Édition manuelle d'un ensemble de 100 lettres inédites complétant l'actuel corpus
- ▶ La base de cas est la base actuelle du corpus
- ▶ Pour chaque version du système de suggestions présentées dans un ordre aléatoire et inconnu, l'évaluateur :
  - édite l'ensemble des lettres ;
  - complète un questionnaire à propos du système ;
  - Attribue un score relatif à la pertinence des suggestions pour chaque propriété du système.

# Évaluation humaine

## Méthodologie

- ▶ Édition manuelle d'un ensemble de 100 lettres inédites complétant l'actuel corpus
- ▶ La base de cas est la base actuelle du corpus
- ▶ Pour chaque version du système de suggestions présentées dans un ordre aléatoire et inconnu, l'évaluateur :
  - édite l'ensemble des lettres ;
  - complète un questionnaire à propos du système ;
  - Attribue un score relatif à la pertinence des suggestions pour chaque propriété du système.

## Résultats

	Basique	Déductif	À base de cas	Combiné
Score moyen (De 1 à 7)	3,4	5,7	5,3	7

# Évaluation automatique

## Méthodologie

- ▶ 200 lettres aléatoirement extraites du corpus
- ▶ Les triplets correspondants sont utilisés pour simuler des questions d'édition
- ▶ Le rang de la valeur attendue est sauvegardé pour chaque version du système de suggestions

# Évaluation automatique

## Méthodologie

- ▶ 200 lettres aléatoirement extraites du corpus
- ▶ Les triplets correspondants sont utilisés pour simuler des questions d'édition
- ▶ Le rang de la valeur attendue est sauvegardé pour chaque version du système de suggestions

## Résultats

	Basique	Déductif	À base de cas	Combiné
$\text{rang} \leq 5$	2,7%	19,23%	33,6%	34,7%
$\text{rang} \leq 10$	7,1%	21,15%	43,2%	43,2%
$\text{rang} \leq 15$	11,3%	22,11%	49,0%	49,5%

Merci de votre attention !

`nicolas.lasolle@univ-lorraine.fr`