

# Da Triple ad Ontologia

Davide Camino

Giugno 2024

## Abstract

In questi appunti mostriamo alcuni esperimenti fatti in merito alla traduzione di triple del tipo (*soggetto, verbo, complemento*) in un'ontologia. Lo scopo di questa traduzione è strutturare i dati in modo formale per poter poi manipolarli attraverso strumenti sofisticati, sia sviluppati ad hoc da noi sia sviluppati da altri come i potenti reasoner che fanno inferenza sulle ontologie.

## 1 Dati di partenza

Per tradurre i dati grezzi utilizzeremo CDuce, un linguaggio di programmazione sviluppato appositamente per manipolare documenti XML. Per questa ragione le triple che andremo a tradurre dovranno essere salvate in formato XML, presentiamo un breve esempio del documento che andremo a tradurre per esemplificarne la struttura:

Listing 1: Dati Grezzi

```
1 <?xml version="1.0"?>
2 <information>
3   <info>
4     <subject>Leonardo da Vinci</subject>
5     <verb>painting</verb>
6     <complement>the Mona Lisa</complement>
7     <sex>M</sex>
8   </info>
9   <info>
10    <subject>Jane Austen</subject>
11    <verb>wrote</verb>
12    <complement>Pride and Prejudice</complement>
13    <sex>F</sex>
14  </info>
15 </information>
```

Da questo codice possiamo vedere che i dati di partenza costituiscono semplicemente un set di informazioni o asserzioni senza alcuna struttura.

Notiamo inoltre l'aggiunta di un attributo che specifica il sesso della persona, facciamo questo perché uno dei nostri obiettivi è quello di valutare i bias di rappresentazione, siamo quindi interessati a verificare se i soggetti di un certo sesso vengono associati a certi verbi più o meno frequentemente dei soggetti del sesso opposto.

## 2 Ontologia di arrivo

Descriviamo brevemente come sarà fatta l'ontologia che vorremmo ottenere.

### 2.1 Struttura

L'ontologia che vogliamo creare per dare una struttura formale ai dati grezzi sarà costituita, come ogni ontologia da classi, individui e relazioni.

### 2.1.1 Classi

L'ontologia che vogliamo ottenere deve rappresentare 2 classi fondamentali:

- **Persone:** la classe che contiene i soggetti delle nostre triple, rappresenta il concetto di persona e tutte le caratteristiche che derivano dall'appartenere alla specie *Homo Sapiens Sapiens*;
- **Complementi:** questa classe rappresenta gli oggetti su cui sono svolte le azioni. è una classe estremamente eterogenea dato che non conterrà solamente oggetti fisici, ma anche astratti quali idee, concetti e pensieri

### 2.1.2 Individui

Gli individui apparterranno a una delle due classi definite precedentemente, vogliamo fare in modo che se due triple di partenza parlano dello stesso soggetto o dello stesso oggetto questo sia rappresentato da un solo individuo anche nell'ontologia che creiamo. Ad esempio:

- Crick, Franklin e Watson scoprirono la struttura del DNA, da queste 3 affermazioni vorremmo ottenere 3 individui di classe Persona e un solo individuo di classe Complemento;
- Albert Einstein scoprì l'effetto fotoelettrico e ideò la teoria della relatività, da queste 2 triple vorremmo ricavare 2 complementi ma un solo soggetto.

### 2.1.3 Relazioni

Per ogni tripla vogliamo aggiungere una relazione alla nostra ontologia che colleghi il soggetto al complemento mediante la specifica azione; Ovviamente vorremmo collegare più soggetti allo stesso complemento (anche mediante azioni differenti) e collegare più complementi allo stesso soggetto.

## 2.2 Formalismo

Esistono diversi modi per esprimere struttura e contenuto di un'ontologia, alcuni formalismi sono più espressivi di altri e pagano questo aumento di espressività con una complessità maggiore sia per quanto riguarda la struttura stessa dell'ontologia che per ciò che riguarda gli algoritmi che faranno inferenza e/o reasoning sui dati contenuti nell'ontologia.

Come già detto nella sezione precedente il linguaggio di programmazione che usiamo è specifico per la manipolazione di documenti XML; prendiamo quindi spunto dalle ontologie create con protégé in particolare valuteremo pro e contro di documenti creati attenendosi ai formati di protégé RDF/XML e OWL/XML

## 3 Linguaggio

Il linguaggio di programmazione scelto per la traduzione è CDuce, un linguaggio di programmazione funzionale general purpose sviluppato appositamente per manipolare documenti XML.

CDuce è in grado di fare il parsing della struttura ad albero di un documento XML in modo molto intuitivo, ogni elemento XML avrà la forma `< (tag)(attr) > content` dove *tag*, *attr* e *content* sono espressioni del linguaggio, in particolare essendo CDuce fortemente tipato qualsiasi elemento XML è del tipo `<(Atom) ({..})>[Any*]`.

Riportiamo qui alcune caratteristiche dei tipi citati, perché sono importanti per comprendere le scelte di progettazioni future

- **Atom:** sono elementi simbolici, vengono usati per specificare i nomi dei tag XML e seguono le stesse regole per gli identificatori delle variabili
- **Record (..):** sono set finiti di (nome, valore) dove i nomi sono label che seguono le stesse convenzioni degli identificatori e i valori sono espressioni
- **Sequences [Any\*]:** sono un elemento fondamentale della sintassi di CDuce, servono per rappresentare le stringhe di caratteri e il contenuto degli elementi XML. Nonostante la centralità delle liste queste sono solo zucchero sintattico per il tipo `Pair` (coppie)

Usiamo questo linguaggio per le garanzie offerte in fase di compilazione in modo da essere sicuri che le funzioni che scriviamo per trasformare elementi ci permettano di ottenere esattamente il tipo di dato desiderato.

CDuce ci permette di descrivere con precisione la struttura del documento XML di partenza e quella del documento di arrivo, le funzioni che mapperanno gli input negli output saranno controllate in fase di compilazione e questo ci assicura che la traduzione avvenga coerentemente con i tipi che abbiamo usato per descrivere i documenti.

## 4 Traduzione

Presentiamo il codice per tradurre le triple di partenza nei due formati presentati. Discutiamo delle implicazioni delle due traduzioni in seguito

### 4.1 RDF/XML

Listing 2: RDF/XML persona

```

1 let fun info_to_pers (Info -> Individual)
2   <info>[ <subject> sub <verb> v <complement> comp <sex> s ] ->
3     let ab = base @ "#" @ sub in
4     let class_res = base @ "#Person" in
5     let prop_res = base @ "#" @ comp in
6     let prop = match v with
7       | "painted" -> <ontology:painted rdf:resource=prop_res> []
8       | "wrote" -> <ontology:wrote rdf:resource=prop_res> []
9       | "discovered" -> <ontology:discovered rdf:resource=prop_res> []
10      | "invented" -> <ontology:invented rdf:resource=prop_res> []
11      | "married" -> <ontology:married rdf:resource=prop_res> []
12      | "built" -> <ontology:built rdf:resource=prop_res> []
13      | _ -> <ontology:action rdf:resource=prop_res> []
14   in
15   <owl:NamedIndividual rdf:about=ab> [
16     <rdf:type rdf:resource=class_res> []
17     prop
18     <rdfs:comment>s
19   ]

```

Applicato questa funzione alla prima informazione del Listato 1 otteniamo:

Listing 3: Traduzione persona RDF

```

1 <owl:NamedIndividual
2   ↪ rdf:about="http://www.semanticweb.org/ontology#Leonardo_da_Vinci">
3   <rdf:type rdf:resource="http://www.semanticweb.org/ontology#Person"/>
4   <ontology:painted
5     ↪ rdf:resource="http://www.semanticweb.org/ontology#the_Mona_Lisa"/>
6   <rdfs:comment>M</rdfs:comment>
7 </owl:NamedIndividual>

```

Con questo codice trasformiamo un elemento di tipo `Info` in un elemento di tipo `Individual`, cioè passiamo da una tripla a un individuo della classe persone. Per completare la traduzione abbiamo poi bisogno di un'altra funzione per trasformare i complementi in individui.

Vediamo che in questo formalismo descriviamo l'azione che collega il soggetto al complemento nello stesso elemento XML che descrive anche la persona. Questo permette una sintassi più compatta e leggibile, ma ci costringe a conoscere a prescindere tutti i verbi che possono collegare un soggetto a un complemento (elencati da riga 7 a 13 del Listato 4), tutti i verbi non contemplati possono essere mappati in un solo verbo, perdendo però la specificità dell'azione.

La ragione per cui si devono conoscere i verbi è dovuta al fatto che in questo formalismo il tipo di azione è rappresentato da tag, come abbiamo accennato questi elementi sono di tipo `Atom`, gli elementi di questo tipo rappresentano la struttura ad albero del documento XML e non possono essere creati durante l'esecuzione del programma perché non si potrebbe garantire la correttezza della trasformazione a nel momento della compilazione.

## 4.2 OWL/XML

Listing 4: RDF/XML persona

```
1 let fun info_to_pers (Info -> Individual)
2   <info>[ <subject> sub <verb> v <complement> comp <sex> s ] ->
3     let ab = base @ "#" @ sub in
4     let class_res = base @ "#Person" in
5     let prop_res = base @ "#" @ comp in
6     let prop = match v with
7       | "painted" -> <ontology:painted rdf:resource=prop_res> []
8       | "wrote" -> <ontology:wrote rdf:resource=prop_res> []
9       | "discovered" -> <ontology:discovered rdf:resource=prop_res> []
10      | "invented" -> <ontology:invented rdf:resource=prop_res> []
11      | "married" -> <ontology:married rdf:resource=prop_res> []
12      | "built" -> <ontology:built rdf:resource=prop_res> []
13      | _ -> <ontology:action rdf:resource=prop_res> []
14    in
15    <owl:NamedIndividual rdf:about=ab> [
16      <rdf:type rdf:resource=class_res> []
17      prop
18      <rdfs:comment>s
19    ]
```

Applicato questa funzione alla prima informazione del Listato 1 otteniamo:

Listing 5: Traduzione persona RDF

```
1 <owl:NamedIndividual
2   ↪ rdf:about="http://www.semanticweb.org/ontology#Leonardo_da_Vinci">
3   <rdf:type rdf:resource="http://www.semanticweb.org/ontology#Person"/>
4   <ontology:painted
5     ↪ rdf:resource="http://www.semanticweb.org/ontology#the_Mona_Lisa"/>
6   <rdfs:comment>M</rdfs:comment>
7 </owl:NamedIndividual>
```

Come si vede in questo caso separiamo la descrizione della persona, da quella dell'azione da quella del complemento. In questo formalismo le azioni non devono essere note a priori, essendo l'azione il valore di un elemento del record e non più un tag.

Come si vede dalla traduzione questo formalismo disperde di più le informazioni e risulta di conseguenza più verboso e meno comprensibile, il vantaggio, come accennato prima è quello di poter descrivere una qualsiasi azione con precisione senza dover conoscere tutte le azioni a priori. Questo formalismo risulta più espressivo e generale; offre la possibilità di descrivere manipolazioni in modo più flessibile con CDuce continuando comunque a poter beneficiare del controllo statico sui tipi offerto dal linguaggio.

## 4.3 Differenze

Possiamo individuare la differenza principale tra i due formati nella gestione delle relazioni tra soggetto e complemento:

- **RDF**: in questo formato è il soggetto l'elemento principale della nostra descrizione, tra i vari elementi che descrivono il soggetto abbiamo i campi che riguardano le relazioni a cui il soggetto partecipa, in questo senso è il soggetto che è collegato tramite un'azione al complemento. Questo approccio è dovuto al fatto che RDF è un formalismo più di basso livello in cui le relazioni non sono la parte centrale della descrizioni degli individui
- **OWL**: questo formato invece pone enfasi sulle relazioni che si instaurano tra gli individui descritti nell'ontologia; soggetto e complemento sono individui separati e indipendenti. L'azione che in questo formato è l'elemento principale della descrizione ed è collegata al soggetto e al complemento in modo da unirli.

## 5 Manipolazione

Ora che abbiamo ottenuto delle ontologie dalle informazioni partenza possiamo fare in modo da sfruttare la struttura formale per poter fare inferenza sui dati, manipolarli e cercare informazioni in modo più efficace rispetto al set di partenza.

In questa sezione mostriamo come implementare una semplice query per cercare tutti i pittori nella nostra ontologia, sfruttiamo questo esempio molto semplice per confrontare 3 differenti approcci, che sfruttino l'ontologia appena creata.

### 5.1 Description Logic

Dato che abbiamo ottenuto un'ontologia (a prescindere che sia in formato RDF o OWL) possiamo adesso sfruttare strumenti molto potenti sviluppati appositamente per lavorare su queste rappresentazioni formali dei concetti. In questo esempio sfruttiamo Protegé e uno dei reasoner offerti da quest'ultimo (HermiT). Grazie al reasoner possiamo interrogare l'ontologia attraverso delle query in Description Logic (DL). Per trovare tutti i pittori ad esempio possiamo semplicemente chiedere:

`painted some`

e otterremo come risposta tutti gli individui che partecipano alla relazione "painted"

### 5.2 CDuce & RDF

Vediamo ora come sia possibile ottenere le stesse informazioni utilizzando partendo da un'ontologia in formato RDF attraverso una funzione in CDuce

Questo codice ci permette di ottenere la lista delle persone che partecipano alla relazione "painted". In questo caso usiamo le query offerte da CDuce per trovare tutti gli individui che facciano il match con l'espressione descritta nella clausola where della query. Come si vede il codice è piuttosto semplice e si basa sul fatto che conosciamo esattamente il tipo del tag della relazione.

In questo formato le azioni sono ristrette a quelle esattamente descritte nella funzione di trasformazione, questa rigidità nella creazione ci permette una relativa snellezza nella ricerca delle informazioni; i problemi sorgono nel caso in cui l'ontologia che vogliamo manipolare non è stata creata da noi ma ci viene assegnata, in questa situazione dobbiamo conoscere esattamente il tipo della relazione cercata per costruire una query ad hoc.

#### 5.2.1 CDuce & OWL

Se il formato dell'ontologia che vogliamo manipolare è OWL possiamo scrivere la funzione di ricerca in questo modo:

## 6 Conclusioni

### 6.1 Risultati ottenuti

In questo lavoro siamo riusciti a arricchire un set di informazioni dando loro una struttura formale grazie alla quale la ricerca e la manipolazione sia supportata da potenti strumenti sviluppati appositamente per le ontologie.

Abbiamo ottenuto differenti formati dell'ontologia, ognuno con vantaggi e svantaggi specifici. Il primo formato descritto è più rigido per quanto riguarda la creazione e la manipolazione, mentre il secondo ci permette più flessibilità.

Dal punto di vista della complessità i formati si scambiano, il primo è più semplice ed è il più adatto ad altre manipolazioni machine oriented, può ad esempio essere l'input di programmi di apprendimento automatico. Il secondo più espressivo risulta più complesso e non particolarmente adatto a algoritmi di AI.

Seppure sia possibile tradurre un formato nell'altro attraverso strumenti differenti da CDuce si vorrebbe, per questioni di overhead evitare una catena continua di traduzioni per passare continuamente da un formato all'altro in base a quello più comodo per la specifica elaborazione.

### 6.2 Lavoro futuro

#### 6.2.1 Duplicati nella traduzione

Per quanto riguarda la traduzione attualmente si fa in modo che due triple differenti rappresentino lo stesso soggetto semplicemente duplicando la rappresentazione, gli strumenti di visualizzazione delle ontologie

sono in grado di evidenziare questi duplicati e eliminarli automaticamente, sarà necessario modificare le funzioni di traduzione in modo che questi duplicati non vengano proprio creati.

### **6.2.2 Struttura dell'ontologia**

Dato che adesso stiamo rappresentando i dati con una struttura formale possiamo aggiungere sottoclassi per raggruppare gli individui in base a particolari caratteristiche di nostro interesse, ad esempio per uno studio sui bias di rappresentazione potrebbe essere utile raggruppare le persone in categorie per valutare se ci sono differenze sostanziali nelle azioni che coinvolgono le varie sottoclassi di persone.