

# Ricerca di nodi vicini

## Knowledge Engineering

**Sofia Galante**

`sofia.galante@stud.unifi.it`

**Agata Parietti**

`agata.parietti@stud.unifi.it`

# Indice

<b>1</b>	<b>Introduzione</b>	<b>3</b>
<b>2</b>	<b>Embeddings</b>	<b>4</b>
2.1	Creazione del dataframe dalla KB . . . . .	4
2.2	TransE . . . . .	5
2.2.1	Implementazione . . . . .	6
2.3	RDF2Vec . . . . .	7
2.3.1	Implementazione . . . . .	7
<b>3</b>	<b>Estensione di SPARQL</b>	<b>9</b>
<b>4</b>	<b>Esempi di query SPARQL con le classi implementate</b>	<b>10</b>
<b>5</b>	<b>Future implementazioni</b>	<b>12</b>
5.1	Estendere a nuovi tipi di Embeddings . . . . .	12
5.2	Ottimizzazione dell'uso della memoria . . . . .	12

# 1 Introduzione

Questo progetto ha come scopo quello di implementare una ricerca di nodi vicini in SPARQL data una Knowledge Base a propria scelta.

Per poter raggiungere questo scopo si deve:

1. eseguire l'embedding dei nodi della KB;
2. estendere SPARQL con una funzione esterna per poter calcolare la distanza tra gli embeddings dei nodi.

Visto che l'embedding svolto mappa ogni nodo in un vettore in uno spazio  $n$ -dimensionale, si è deciso di utilizzare la **cosine distance** come metrica. La formula è la seguente:

$$d(v_1, v_2) = 1 - \cos(\theta) = 1 - \frac{v_1 \cdot v_2}{||v_1|| ||v_2||} \quad (1)$$

dove  $\theta$  è l'angolo tra i vettori  $v_1$  e  $v_2$ .

Si noti che la distanza è quindi un numero reale che sta tra 0 (nodi identici) e 2 (nodi opposti).

Gli embeddings vengono svolti con dei programmi python, mentre l'estensione di SPARQL è stata fatta con delle classi Java. L'estensione è stata fatta per la piattaforma Fuseki.

## 2 Embeddings

Per decidere l'embedding giusto da usare si sono letti diversi articoli (presenti in bibliografia). Durante la ricerca, si è notato che esistono due diverse scuole di pensiero per quanto riguarda gli embeddings di una KB:

1. la prima vede la KB come un Knowledge Graph e cerca di creare un embedding che rispetti le regole di transizione nel grafo;
2. la seconda trasforma le triple della KB in frasi e svolge un embedding classico dell'NLP, basato sui language models.

Per questo motivo, si è deciso di implementare un metodo per entrambe le scuole di pensiero, lasciando poi all'utente la scelta di quale utilizzare durante la sua ricerca.

### 2.1 Creazione del dataframe dalla KB

Per poter effettuare l'embedding sulla KB, si deve prima trasformarla in un dataframe di triple.

Per compiere questo passaggio, si è utilizzata la libreria *rdflib* che permette di leggere un grafo di una KB e di estrarre da essa tutte le triple presenti. Il risultato è poi stato salvato con *pandas*.

Nel nostro progetto, è possibile sia creare un embedding di un grafo salvato in locale (in questo caso il grafo dovrà essere inserito nella cartella *dataset*) oppure di utilizzare un grafo direttamente da un rdfstore.

La seconda modalità è stata resa possibile dall'utilizzo della libreria *sparqlwrapper*. In caso si volesse eseguire l'embedding di un grafo direttamente da un rdfstore, i passaggi sono i seguenti:

1. si deve modificare il file *sparql.txt*, contenente l'rdfstore di interesse;
2. si deve modificare il file *query.txt*, contenente una query scritta in linguaggio SPARQL; si noti che questa query deve essere di tipo CONSTRUCT, in quanto rdflib deve ricevere in ingresso un grafo serializzato nel formato RDF.

Un esempio di file *sparql.txt* e di *query.txt* è il seguente:

1. **sparql.txt:** <https://query.wikidata.org/sparql>

2. **query.txt:**

```
CONSTRUCT {?s ?p ?o}  
WHERE {  
  ?s ?p ?o  
}  
LIMIT 100000}
```

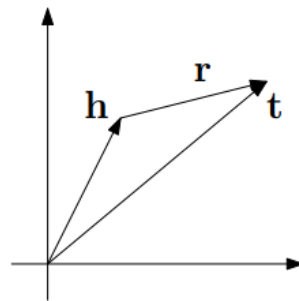
In questo caso, il grafo costruito è composto dalle prime 100000 triple presenti su *wikidata*.

## 2.2 TransE

L'embedding di tipo TransE è un embedding di tipo *transizionale*.

Con l'algoritmo TransE, ogni entità e ogni relazione di una KB viene mappata in uno spazio vettoriale facendo in modo che si mantenga la relazione tra le triple.

Data una tripla  $(h, r, t)$  vengono fatti degli embedding per ogni elemento  $(\mathbf{h}, \mathbf{r}, \mathbf{t})$  tali che:



$$\mathbf{h} + \mathbf{r} \simeq \mathbf{t}$$

Per imparare gli embedding si minimizza un margin-based ranking criterion

$$L = \sum_{(h,r,t) \in S} \sum_{(h',r,t') \in S_{(h,r,t)}} [\gamma + d(\mathbf{h} + \mathbf{r}, \mathbf{t}) - d(\mathbf{h}' + \mathbf{r}, \mathbf{t}')]_+ \quad (2)$$

dove:

1.  $S$  è l'insieme di tutte le triple della KB;
2.  $(\mathbf{h}', \mathbf{r}, \mathbf{t}')$  è una tripla corrotta (cioè ottenuta prendendo una tripla esistente in  $S$  e modificando uno dei due nodi) e  $S_{(h,r,t)}$  è l'insieme di triple corrotte ottenute a partire dalla tripla  $(h, r, t)$ ;
3.  $d(\mathbf{h} + \mathbf{r}, \mathbf{t})$  è una misura di dissimilarità (nel nostro caso la  $L_2 - norm$ );
4.  $\gamma$  è un margin hyperparameter;
5.  $[x]_+$  indica la parte intera di  $x$ .

### 2.2.1 Implementazione

L'algoritmo che implementa TransE si trova nel file *te\_embedding.py*.

Il modello TransE si trova nella libreria *torchkge*, un'estensione di *PyTorch* che permette di fare vari embeddings di tipo transizionale (e non) su una KB.

Gli iperparametri che l'utente può scegliere sono:

1. *emb\_dim*  $\rightarrow$  la dimensione dei vettori dell'embedding;
2. *lr*  $\rightarrow$  learning rate per il training;
3. *n\_epochs*  $\rightarrow$  numero di epoche per il training;
4. *b\_size*  $\rightarrow$  dimensione del batch;
5. *margin*  $\rightarrow$  margin hyperparameter.

Dei valori di default sono stati scelti per la KB di prova fornita da fuseki (**Cheese.ttl**) ma si consiglia di modificare i parametri a seconda delle dimensioni della KB utilizzata.

Si noti che la KB viene prima trasformata in un Dataframe pandas per poter essere data in input al modello. La trasformazione in Dataframe pandas è stata resa possibile dal parser presente nella libreria *rdflib* e l'implementazione di questa trasformazione si trova nel file *dataset\_preparation.py*.

## 2.3 RDF2Vec

RDF2Vec è un'interpretazione per KB dell'embedding NLP Word2Vector.

In questo caso le triple della KB vengono reinterpretate come delle frasi della forma  $nodo_1 - relazione_1 - nodo_2 - relazione_2 - nodo_3$  e così via (si noti che la lunghezza della catena  $nodo - relazione$  è un iperparametro).

L'algoritmo RDF2Vec compie quindi i seguenti passi:

1. per ogni nodo, si creano  $n$  "frasi" di lunghezza  $d$  ( $n$  e  $d$  sono a scelta dell'utente);
2. su queste frasi si applica il Word2Vector  $\rightarrow$  gli embedding vengono aggiornati per imparare un *language model* che permette di prevedere la "parola" (i.e. il nodo o la relazione) che si trova al posto  $t$  di una frase date le parole che ha prima, dopo, o intorno.

Nel nostro caso, si è deciso di utilizzare il *Continuous Bag-of-Words Model*, in cui la parola predetta dal modello è la parola centrale in una finestra di grandezza predefinita.

Si vuole, quindi, massimizzare la *average log-probability*:

$$\frac{1}{T} \sum_{t=1}^T \log P(w_t | w_{t-c}, w_{t-c+1}, \dots, w_{t+c-1}, w_{t+c}) \quad (3)$$

dove  $T$  è il numero di parole totali.

La probabilità sopra riportata è calcolata con una *softmax*.

### 2.3.1 Implementazione

L'algoritmo che implementa RDF2Vec si trova nel file `w2v_embedding.py`.

La libreria utilizzata per compiere l'embedding è *gensim*.

Gli iperparametri che l'utente può scegliere sono:

1.  $n\_phrases \rightarrow$  numero di "frasi" costruite per nodo;
2.  $dim\_phrases \rightarrow$  numero di "parole" presenti in una "frase";
3.  $workers \rightarrow$  numero di thread della CPU usati per fare l'embedding;

4. *emb\_dim* → la dimensione dei vettori dell'embedding;
5. *min\_count* → il numero minimo di volte con cui una "parola" deve apparire nelle "frasi" per essere considerata;
6. *sg* → un parametro che indica il metodo con cui creare il *language model* (0 = Continuous Bag-of-Words Model, 1 = Skip-Gram).

Dei valori di default sono stati scelti per la KB di prova fornita da fuseki (**Cheese.ttl**) ma si consiglia di modificare i parametri a seconda delle dimensioni della KB utilizzata.

Si noti che la KB viene prima trasformata in un Dataframe pandas per poter essere data in input al modello. La trasformazione in Dataframe pandas è stata resa possibile dal parser presente nella libreria *rdflib* e l'implementazione di questa trasformazione si trova nel file *dataset\_preparation.py*.



### 3 Estensione di SPARQL

L'estensione di SPARQL è stata fatta attraverso 3 diverse classi java: *Embeddings.class*, *W2V.class* e *TE.class*.

La classe *Embeddings.class* contiene tutte le funzioni per caricare gli embeddings dei nodi e calcolare la distanza tra di essi, mentre le classi *W2V.class* e *TE.class* sono estensioni della classe *FunctionBase2* di *fuseki* e implementano la funzione *exec(NodeValue1, NodeValue2)* che prende in ingresso due nodi della KB e restituisce un terzo nodo in uscita (nel nostro caso, il nodo in uscita è un numero *double* equivalente alla distanza tra i nodi in ingresso).

Per una gestione ottimale della memoria, si noti che gli embeddings vengono caricati solo al primo utilizzo della funzione e poi salvati in una *HashMap* persistente in memoria.

Per funzionare, *fuseki* deve accedere unicamente ai file *.class* e *.jar* delle classi implementate. Si è comunque deciso di condividere il file sorgente *.java* per poter vedere più nel dettaglio il codice.

Nella sezione successiva sono presenti diversi esempi su come si possono utilizzare le classi.

## 4 Esempi di query SPARQL con le classi implementate

Le classi  $TE$  e  $W2V$  possono essere richiamate all'interno di normali query SPARQL.



```
1 PREFIX f: <java:>
2
3 SELECT DISTINCT ?dist WHERE {
4   BIND(f:W2V(<http://data.kasabi.com/dataset/cheese/halloumi>, "Halloumi"@en) as ?dist)
5 } ORDER BY ?dist
```

Table Response 1 result in 0.041 seconds Simple view Ellipse Filter query results Page size: 50

dist
1 "0.10986871190415715e0"^^<http://www.w3.org/2001/XMLSchema#double>

Come mostrato in figura, la funzione restituisce la distanza tra due nodi.

I risultati ottenibili autonomamente dalla funzione sono alquanto limitati, ma si possono estendere utilizzando le clausole SPARQL e avere risultati molto più completi.

Ad esempio, si può ottenere l'elenco dei nodi vicini ad un nodo ordinati in base a quanto sono effettivamente vicini:

```

1 PREFIX f: <java:>
2
3 SELECT DISTINCT ?s1 ?dist WHERE {
4   ?s1 ?p1 ?o1
5   BIND(f:TE(<http://data.kasabi.com/dataset/cheese/halloumi>, ?s1) as ?dist)
6   FILTER(?s1 != <http://data.kasabi.com/dataset/cheese/halloumi>)
7 } ORDER BY ?dist

```

s1	dist
<http://ec.europa.eu/eurostat/ramon/rdfdata/nuts2008/ITC4A>	"0.1492566724703066e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://ec.europa.eu/eurostat/ramon/rdfdata/nuts2008/ITE42>	"0.1948317100116731e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/formaggella-del-luinese>	"0.20044471778829642e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/italy/pesaro>	"0.22127851360444872e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/italy/pradlevies>	"0.2219525213230934e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://ec.europa.eu/eurostat/ramon/rdfdata/nuts2008/ITE2>	"0.25230508622893744e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/italy/massa>	"0.2649628983587742e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://ec.europa.eu/eurostat/ramon/rdfdata/nuts2008/ITD5>	"0.286418975389101e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/italy/paulliatino>	"0.28714477591737597e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/italy/barletta-andria-trani>	"0.2911424387829763e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://ec.europa.eu/eurostat/ramon/rdfdata/nuts2008/ITD3>	"0.29912183902101497e0"^^<http://www.w3.org/2001/XMLSchema#double>

Oppure si possono calcolare le distanze tra tutti i nodi nella KB e osservare quali sono i nodi più vicini tra loro:

```

1 PREFIX f: <java:>
2 SELECT ?s1 ?s2 ?dist WHERE {
3   {
4     SELECT ?s1 ?s2
5     WHERE {
6       ?s1 ?p1 ?o1.
7       ?s2 ?p2 ?o2
8       FILTER(?s1 != ?s2)
9     } GROUP BY ?s1 ?s2
10    ORDER BY ?s1
11  }
12  BIND(f:W2V(?s1, ?s2) as ?dist)
13 }
14 LIMIT 25

```

s1	s2	dist
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/afuegal-pitu>	"0.10265697236641502e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/agri-di-valtorta>	"0.026966752136078576e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/allgauer-bergkase>	"3.731351638780467E-8"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/allgauer-ementaler>	"0.00835770831278193e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/allgauer-sernalpkase>	"0.043268288673893096e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/almkase>	"0.005817892214388443e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/altenburger-ziegenkase>	"3.675549204948947E-4"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/anevato>	"0.002642369172318526e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/arzua-ulloa>	"0.05775541015736352e0"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/asiago>	"6.136143540659278E-4"^^<http://www.w3.org/2001/XMLSchema#double>
<http://data.kasabi.com/dataset/cheese/abondance>	<http://data.kasabi.com/dataset/cheese/asiago-dallevo>	"0.07569278082016029e0"^^<http://www.w3.org/2001/XMLSchema#double>

Infine, si noti che per passare da una funzione all'altra basta modificare il nome della funzione chiamata all'interno del BIND.

## 5 Future implementazioni

### 5.1 Estendere a nuovi tipi di Embeddings

Il progetto da noi presentato è estendibile sia nei programmi python, sia nelle classi java per essere utilizzato con altri tipi di embeddings diversi da *RDF2Vec* e *TransE*.

In particolare, la libreria *torchkge* utilizzata per *TransE* contiene al suo interno altri modelli di embeddings, rendendo il tutto, quindi, ancora più facile da estendere.

Per quanto riguarda le classi java, l'unica cosa da fare è creare una classe java che estende la classe *FunctionBase2* di Fuseki e che nella sua funzione *exec()* richiami la funzione omonima della classe *Embeddings*.

### 5.2 Ottimizzazione dell'uso della memoria

Nel caso di utilizzo di una KB salvata in locale, i file python attuali caricano l'intera KB per crearne l'embedding. Questo può essere un problema per KB molto grandi, in quanto occupano molta memoria in RAM.

Una possibile implementazione futura può essere, quindi, quella di ottimizzare la lettura della KB in memoria.

## Riferimenti bibliografici

- [1] Apache-Jena Fuseki  
<https://jena.apache.org/documentation/fuseki2/index.html>
- [2] Petar Ristoski, Heiko Paulheim. RDF2Vec: RDF Graph Embeddings for Data Mining  
[https://madoc.bib.uni-mannheim.de/41307/1/Ristoski\\_RDF2Vec.pdf](https://madoc.bib.uni-mannheim.de/41307/1/Ristoski_RDF2Vec.pdf)
- [3] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran. Translating Embeddings for Modeling Multi-relational Data  
<https://proceedings.neurips.cc/paper/2013/file/1cecc7a77928ca8133fa24680a88d2f9-Paper.pdf>
- [4] TorchKGE  
<https://torchkge.readthedocs.io/en/latest/>
- [5] Gensim  
<https://radimrehurek.com/gensim/>
- [6] RDFLib  
<https://rdflib.readthedocs.io/en/stable/>
- [7] SparqlWrapper  
<https://pypi.org/project/SPARQLWrapper/>