

TLN TerzaParte

Mario Scapellato

Terza Parte dell'esame in Tecnologia del Linguaggio Naturale.
Prof. Luigi Di Caro

1 Prima e Seconda Esercitazione

1.1 Introduzione

La prima e la seconda esercitazione riguardano prima di tutto la creazione di definizione per 4 concetti scelti da noi studenti a lezione, dando, per ciascun concetto, 2 asserzioni concrete (di cui 1 generico e 1 specifico) e 2 astratte (di cui 1 generico e 1 specifico). Ogni studente ha inserito all'interno del documento una definizione per ciascun concetto in maniera indipendenten e sulla base della propria conoscenza sull'argomento stesso.

I 4 concetti scelti sono i seguenti :

- Generico Astratto : *Courage*
- Generico Concreto : *Paper*
- Specifico Astratto : *Apprehension*
- Specifico Concreto : *Sharpener*

Nel file defs.csv sono presenti le 30 definizioni assegnate per i 4 concetti illustrati.

1.2 Sviluppo

L'esercitazione e' stata sviluppata nella seguente maniera :

- Sono partito prima di tutto con una semplice implementazione di una funzione che mi permettesse di leggere riga per riga il mio csv.
- Successivamente ho eseguito le varie operazioni di Text Cleaning, ovvero :
 - *remove_punctuation* : operazione che mi permette di rimuovere tutte le punteggiature; data una frase contenente le punteggiature, verra' restituita una senza punteggiatura.
 - *remove_stopword* : mediante il file stop_words_FULLL.txt ho eliminato tutte le stopwords per ogni frase data in input.

- *tokenize_sentences*: suddivido una frase in unita' piu' piccole (come singole parole) e sucessivamente, ogni parola viene portata al suo lemma
- Sucessivamente ho definito la funzione *definitions* che mi legge il file .csv, definendo un dizionario formato da una coppia chiave-valore, dove la chiave e' il termine, mentre il valore rappresenta una lista di *Bag Of Words* correlate alle definizioni.
- Ho calcolato poi la similarita' del coseno : essa e' una misura che viene utilizzata per misurare la similarita' tra documenti durante la fase di text analysis. La formula e' la seguente :

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}},$$

La similarita' del coseno prende in input 2 vettori inizialmente vuoti che corrispondono alle due definizioni che poi verranno prese in analisi. I vettori conterranno valore 1 se la parola analizzata appartiene al set, 0 altrimenti. Questo procedimento viene svolto per ogni singola parola che appare in entrambe le definizioni. Una volta fatta questa analisi, posso andare a calcolare della distanza del coseno e restituirne i risultati.

- Nella funzione definita *compute_results*, prendo in input il dizionario di liste contenente tutte le definizioni processate e vado a calcolare la similarita' del coseno tra tutte le coppie di definizioni dello stesso concetto. Verra' restituito il valore medio della similarita' per ogni concetto.
- Nella funzione *most_frequent_words*, prendo in considerazione l'insieme delle parole presenti all'interno di ciascuna definizione e considero in input un dizionario che conterra' una lista di parole frequenti per tutte le definizioni. Per ogni concetto, vado a calcolare quelle parole che sono presenti in almeno il 50% delle definizioni. Verra' restituito l'elenco delle parole piu' frequenti.

1.3 Risultati Ottenuti

I risultati ottenuti, come annunciato precedentemente a lezioni, non sono particolarmente elevati; in particolare :

```
Similarita del coseno : {'Courage': 0.21054727554969985, 'Paper':
0.29258850377799267, 'Apprehension': 0.0830330313557733, 'Sharpener':
0.3863878711824424}
```

i risultati mostrano di come concetti piu' concreti come *Paper* e *Sharpener* sono descritti mediante delle definizioni piu' semplici e in linea tra loro, a differenza

dei concetti piu' generici come *Apprehension* e *Courage* che trovano maggiori difficolta' ad essere descritti da delle definizioni simili tra di loro. Inoltre, vediamo che le parole piu' frequenti presentano dei risultati piu' "attesi" :

```
Most frequent words :
[('Courage', ['ability', 'fear']), ('Paper', ['write', 'material']),
 ('Apprehension', []), ('Sharpener', ['pencil', 'sharpen', 'tool'])]
```

2 Terza Esercitazione

2.1 Introduzione

Tra le tre proposte, e' stato scelto il task che riguarda la caratterizzazione delle risorse attraverso WordNet. L'obiettivo consiste nel ricercare i vari pattern (che fanno riferimento alle definizioni) all'interno delle ontologie e definire uno studio che riguardo la forma (come la lunghezza), il contenuto o l'aspetto relazionale (come la presenza di iperonimi, antonimi ecc..).

2.2 Sviluppo

La verifica dei task di analisi e' stata fatta mediante l'utilizzo della risorsa lessicale Wordnet.

Gli studi effettuati sono stati principalmente 2:

- Nella prima parte, e' stata definita la lunghezza di ciascuna definizione. L'idea era la seguente : preso in considerazione un concetto (in questo caso erano i concetti che avevamo analizzato nell'esercitazione precedente, ovvero *Courage*, *Paper*, *Sharpner* e *Apprehension*), tramite il metodo *definition_lenght(word)* verificavo la lunghezza di ciascuna definizione presente in WordNet associata a quel concetto.
- Nella seconda parte dell'esercitazione ho focalizzato l'attenzione sugli aspetti relazionali: per ogni concetto, ho osservato il numero di definizioni associate, il numero totale di iperonimi, meronimi e antonimi. Sucessivamente, ho calcolare una piccola percentuale che rappresenta il numero di meronimi,antonimi o iperonimi trovati diviso il numero totale delle definizioni associati al concetto di riferimento.
- Infine ho implementato un metodo aggiuntivo, definito *total_hyponym_path(word)* in cui data una parola, viene considerato ogni suo significato e per ogni suo significato viene calcolato e stampato il suo iperonimo.

2.3 Risultati Ottenuti

I risultati ottenuti sono suddivisi sulla base delle due parti svolte durante l'esercitazione:

```

-----
Concept:  Courage
{15: 1}

-----

Concept:  Paper
{3: 2, 5: 1, 6: 1, 8: 1, 11: 2, 14: 1, 15: 1}

-----

Concept:  Apprehension
{2: 1, 4: 1, 7: 1, 8: 1}

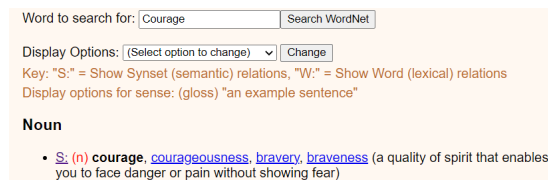
-----

Concept:  Sharpener
{14: 1}

```

Figure 1: Risultati ottenuti primo studio.

Come possiamo vedere dalla Figura 1, per ogni concetto che analizziamo, ne riportiamo la rispettiva lunghezza della definizione (presa da WordNet) associata al concetto con l'ulteriore aggiunta del tag di riferimento associata alla definizione del concetto analizzato. Vediamo di spiegare meglio l'idea: per il concetto *Courage*, vediamo che in forma di dizionario viene stampato il risultato costituito da una coppia chiave-valore 15:1. Se vado su WordNet, e cerco il significato per la parola *Courage*, avrò una definizione lunga 15 parole, e 1 che fa riferimento al tag associato al concetto (in questo caso *Noun*). Lo stesso discorso vale anche per gli altri concetti.



Per quanto riguarda il secondo studio, i risultati hanno mostrato un evidente discrepanza fra il numero di possibili iperonimi, meronimi e antonimi per ciascuna delle definizioni di ciascun senso. Vediamo che abbiamo un evidente predominanza degli iperonimi che risultano essere nettamente maggiori rispetto ai meronimi, mentre per quanto riguarda gli antonimi non vi è nessuna presenza.

```

Concept:  Courage

tot definizioni :1      iperonimi trovati : 1      ==> 1.0
tot definizioni : 1      meronimi trovati : 0      ==> 0.0
tot definizioni : 1      antonimi trovati : 0      ==> 0.0

-----

Concept:  Paper

tot definizioni :9      iperonimi trovati : 7      ==> 0.7777777777777778
tot definizioni : 9      meronimi trovati : 1      ==> 0.1111111111111111
tot definizioni : 9      antonimi trovati : 0      ==> 0.0

-----

Concept:  Apprehension

tot definizioni :4      iperonimi trovati : 1      ==> 0.25
tot definizioni : 4      meronimi trovati : 0      ==> 0.0
tot definizioni : 4      antonimi trovati : 0      ==> 0.0

-----

Concept:  Sharpener

tot definizioni :1      iperonimi trovati : 1      ==> 1.0
tot definizioni : 1      meronimi trovati : 0      ==> 0.0
tot definizioni : 1      antonimi trovati : 0      ==> 0.0

```

Ecco che in figura notiamo che per il primo concetto (*Courage*), abbiamo per una definizione, la presenza di un singolo iperonimo; per *Paper* la situazione e' diversa: vediamo che per un totale di 9 definizioni per ciascun senso, abbiamo 7 iperonimi, 1 meronimi e un'altro che non fa un altro riferimento, avendo un totale di circa il 77% di iperonimi a fronte dell'11% di meronimi. Stesso discorso vale per *Apprehension*: su 4 definizioni trovate, abbiamo 1 solo riferimento per gli iperonimi (quindi il 25%), mentre i restanti 3 non rappresentano ne' meronimi, ne' antonimi. Il concetto *Sharpener* ha lo stesso comportamento di *Courage*, restituendone gli stessi risultati.

Proprio perche' gli iperonimi assumono un ruolo particolarmente importante per l'aspetto relazionale, vediamo quali sono gli iperonimi totali per ciascun concetto. Nella figura che segue verranno riportati per semplicita' solamente quelli per *Courage* :

```

concept: courage
[[{"synset": "entity.n.01", 1}, {"synset": "abstraction.n.00", 1}, {"synset": "attribute.n.00", 0}, {"synset": "trait.n.00", 7}, {"synset": "character.n.01", 10}, {"synset": "spirit.n.00", 9}, {"synset": "courage.n.01", 10}]

```

Se andiamo a dare un'occhiata su WordNet vediamo che i risultati coincidono esattamente :

- [S; \(n\) courage, courageousness, bravery, braveness](#) (a quality of spirit that enables you to face danger or pain without showing fear)
 - [direct hyponym](#) / [full hyponym](#)
 - [attribute](#)
 - [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
 - [S; \(n\) spirit](#) (a fundamental emotional and activating principle determining one's character)
 - [S; \(n\) character, fiber, fibre](#) (the inherent complex of attributes that determines a person's moral and ethical actions and reactions)
 - "education has for its object the formation of character". *Herbert Spencer*
 - [S; \(n\) trait](#) (a distinguishing feature of your personal nature)
 - [S; \(n\) attribute](#) (an abstraction belonging to or characteristic of an entity)
 - [S; \(n\) abstraction, abstract entity](#) (a general concept formed by extracting common features from specific examples)
 - [S; \(n\) entity](#) (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

3 Quarta Esercitazione

3.1 Introduzione

La quarta esercitazione prevede lo studio di alcuni verbi dal punto di vista della Teoria di Hansk. Secondo Hansk il verbo è la radice del significato e non esistono espressioni senza verbo. Ad ogni verbo viene associata una valenza che indica gli argomenti che sono necessari per il verbo. Possiamo differenziare il significato del verbo in base al numero di argomenti che possiede. Determinato il numero di argomenti per un certo verbo, bisogna specificarli mediante un certo numero di slot; ogni slot è formato da un certo numero di valori che lo riempiono, detti *filler*. Ogni filler può avere associati dei tipi semantici che rappresentano delle generalizzazioni strutturate come una gerarchia.

3.2 Sviluppo

Per lo svolgimento sono stati utilizzati i seguenti verbi :

- Eat
- Read

Per recuperare le n frasi in cui il verbo viene utilizzato, è stato utilizzato il corpus brown. Poiché la forma di una parola può cambiare in base al verbo, è necessario l'utilizzo del lemmatizzatore di WordNet. Iterando su tutte le parole di una singola frase, per ogni parola viene calcolato il lemma, e se è uguale alla forma base del verbo scelto, la frase viene salvata.

Per ogni frase, mediante l'ausilio di *spacy*, effettuando il pos e il parsing, ottengo i vari token collegati tramite le dipendenze sintattiche. Iterando sui vari token, è possibile estrarre il soggetto e il complemento oggetto relativi al verbo di una frase. Soggetto e complemento oggetto vengono poi mappati con i relativi super sensi di WordNet nel seguente modo :

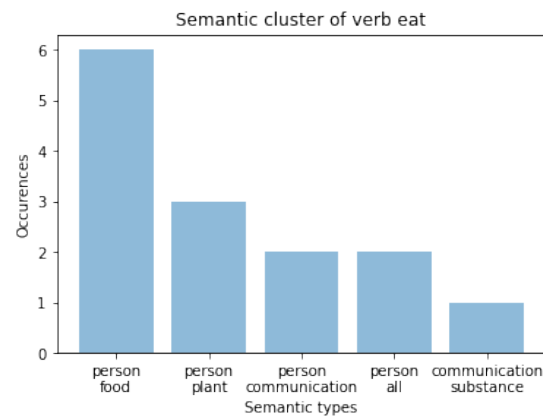
- Nome proprio : Person
- Pronome : ['i', 'you', 'he', 'she', 'we', 'they', 'me']:person, oppure [it]:entity

- Parola generica: in questo caso entra in gioco il WSD che tenta di restituirne il senso del termine data la frase e il termine generico stesso.

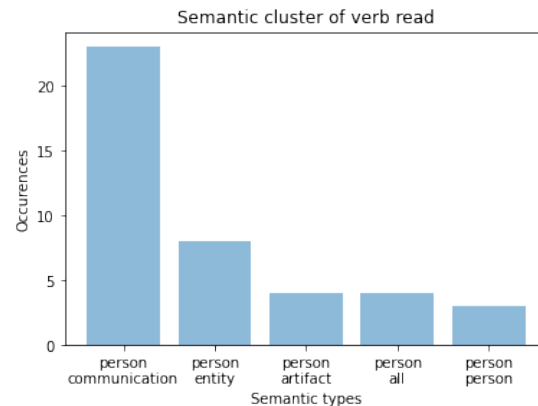
I supersensi di soggetto e complemento oggetto rappresentano i filler dei verbi. Le coppie dei supersensi sono contenute all'interno di cluster semantici opportunamente rappresentati. Successivamente, vengono calcolati per entrambi i filler i supersensi più frequenti.

3.3 Risultati Ottenuti

Possiamo vedere i plot relativi sia per il verbo Eat che per il verbo Read:



chiaramente possiamo notare che la coppia più frequente è quella definita dalla coppia person:food con una frequenza pari a 7.



Per il verbo Read invece la coppia più frequente è quella relativa alla coppia person:communication, con una frequenza pari a 23.

Vediamo quindi i risultati sulle frequenze stampati a video :

```
< communication:substance > Count 1
< communication:quantity > Count 1
< quantity:body > Count 1
< person:plant > Count 3
< person:person > Count 1
< person:animal > Count 1
< group:plant > Count 1
< person:communication > Count 2
< entity:body > Count 1
< person:quantity > Count 1
< person:all > Count 2
< person:food > Count 6
< person:entity > Count 1
< animal:group > Count 1
< person:substance > Count 1
```

```
< person:communication > Count 23
< person:artifact > Count 4
< person:entity > Count 8
< group:communication > Count 1
< group:location > Count 1
< person:change > Count 1
< person:person > Count 3
< entity:person > Count 1
< person:cognition > Count 1
< person:attribute > Count 1
< all:quantity > Count 2
< person:time > Count 1
< person:all > Count 4
< quantity:possession > Count 1
< act:communication > Count 2
< quantity:entity > Count 2
< person:quantity > Count 1
< relation:time > Count 1
< group:state > Count 2
< group:stative > Count 1
< person:substance > Count 1
< person:event > Count 1
< person:location > Count 1
< group:person > Count 1
```


4 Quinta Esercitazione

4.1 Introduzione

La quinta esercitazione prevede la sperimentazione del content-to-form, cioè cercare di risalire al synset di un concetto indirizzando la ricerca in WordNet attraverso i genus. Alla base, vi è il principio *Genus-Differentia definition*, secondo il quale un concetto può essere definito da due elementi principali :

- *Genus* : parte di una definizione esistente che viene utilizzata come porzione per una nuova definizione; tutte le definizioni con lo stesso genus vengono raggruppate secondo lo stesso genus
- *Differentia*: porzione della definizione che non viene dal genus e che rende più specifico un concetto.

4.2 Sviluppo

Sono state implementate le seguenti operazioni :

- Operazioni di Text Cleaning e preprocessing grazie al solito svoglimento di rimozione delle stopwords, rimozione della punteggiatura e tokenizzazione delle varie frasi.
- Una volta che sono state ottenute tutte le definizioni dei concetti, grazie al metodo *get_definitions(file)* svolto nell'esercitazione precedente, posso applicare il metodo di ricerca del genus. Prima di tutto, utilizzo il modulo *Counter* per ottenere un dizionario in cui le chiavi sono le parole che appartengono alla definizione e i valori rappresentano la frequenza di comparsa di quella parola. Impostiamo il valore numerico 5 come il numero di termini più importanti da estrarre per un concetto. La funzione mi restituirà una lista ordinata dei 5 termini più importanti (ovvero quelli più frequenti)
- Per ogni genus, ottengo un possibile candidato; per ogni genere ottengo una lista di iponimi di tutti i synset relativi al termine genus. Verrà restituito quel synset che si avvicina di più alle definizioni del concetto; questo approccio permette di massimizzare l'overlap tra la signature dell'iperonimo e il BoW delle definizioni associate al concetto in analisi

4.3 Risultati Ottenuti

I risultati ottenuti dal task sono i seguenti :

```
concept: courage
genus list (with frequency):
[('ability', 10), ('fear', 17), ('face', 9), ('situation', 7), ('scar', 5)]

candidates:
[('ability', synset('physical_ability.n.01')), ('fear', synset('stage_fright.n.01')), ('face', synset('take_the_bull_by_the_horns.v.01')), ('situation', synset('crowding.n.01')), ('scar', synset('bandaid.n.01'))]

concept: paper
genus list (with frequency):
[('material', 22), ('write', 10), ('collabor', 7), ('wood', 6), ('tree', 5)]

candidates:
[('material', synset('composite_material.n.01')), ('write', synset('handwrite.v.01')), ('collabor', synset('help.n.01')), ('wood', synset('balsa.n.01')), ('tree', synset('palm.n.01'))]
```

```

concept: Apprehension
genus: list (with frequency)
[('fear', 80), ('anxiety', 80), ('feeling', 5), ('happen', 5), ('feel', 4)]

candidates:
[('fear', synset('apprehension.n.01')), ('anxiety', synset('anxiety.n.01')), ('feeling', synset('feeling.v.01')), ('happen', synset('happen.v.01')), ('feel',
synset('feel.v.01'))]
.....
concept: Sharper
genus: list (with frequency)
[('penicil', 25), ('sharper', 17), ('tool', 16), ('object', 11), ('allow', 4)]

candidates:
[('penicil', synset('penicillin.n.01')), ('sharper', synset('sharper.v.01')), ('tool', synset('tool.n.01')), ('object', synset('object.n.01')), ('allow',
synset('allow.v.01'))]

```

Come possiamo vedere, solo in un unico caso troviamo una corrispondenza esatta tra il miglior senso presente nella lista genus e il senso del concetto in analisi. Parliamo proprio del concetto *Apprehension* e il suo genus *fear*: entrambi sono mappati al *Synset('apprehension.n.01')*). Questo può essere dovuto da un legame semantico debole che c'è tra il concetto preso in analisi e l'intorno ottenuto attraverso i genus e i loro iperonimi. Inoltre, il documento *defs.csv* è un file che è stato rappresentato da noi studenti, quindi alcune definizioni possono essere lontane da ciò che WordNet intende per quel determinato concetto.

5 Sesta Esercitazione

La sesta esercitazione prevedeva la realizzazione di un semplice algoritmo di summarization. Il task consiste nell'effettuare un riassunto automatico a partire da un testo in input. Il riassunto viene creato estraendo parti di testo rilevanti (interi paragrafi nel testo) in base a un punteggio che viene assegnato a ciascun paragrafo a seconda delle parole rilevanti che ci sono all'interno di ciascun paragrafo. Il risultato ottenuto offre una compressione del testo del 40%.

5.1 Sviluppo

Per l'implementazione del seguente task ho deciso di riportare in formato .txt una pagina di Wikipedia che trattava di Machine Learning. Successivamente sono state effettuate le seguenti implementazioni :

- Dato il testo, effettuo le solite operazioni di rimozione delle stopword, punteggiatura, oltre grazie al supporto fornito dalla libreria spacy che permette di effettuare le opportune operazioni di estrazione ed elaborazione dei testi analizzati.
- Successivamente, posso passare all'operazione di lettura del mio documento. Per fare ciò utilizzo il modulo glob che mi permette di recuperare file o percorsi corrispondenti ad un modello specificato.
- Una volta aperto il mio documento, conto le parole più frequenti all'interno del mio .txt: la frequenza delle parole mi sarà poi utile per fare la summarization finale.
- Per ciascun elenco di parole, vado a dividere ogni parola per il massimo valore di frequenza ottenuto. Questi punteggi espressi per ogni parola mi servono poi per capire se una determinata parola è presente o no in un

paragrafo e, se e' presente, probabilmente quel paragrafo avra' un rilevanza maggiore rispetto ad un altro.

- Posso creare dei punteggi per ogni frase: il punteggio fa riferimento a quanto l'algoritmo considera rilevante o no quel paragrafo in relazione al contesto su cui sta operando.
- Tramite il modulo *nlargest* prendo i paragrafi con i punteggi migliori, ci applico una compressione del 40% e ne restituisco il riassunto finale
- Applico una metrica di score denominata *rouge_score* per visualizzare i risultati.

5.2 Risultati Ottenuti

Come abbiamo detto in precedenza, il risultato ottenuto e' stato definito a partire da una compressione del testo del 40%. Prima vediamo i punteggi associati ad ogni paragrafo :

```
document { machine learning (ML) is a field of inquiry devoted to understanding and building methods that "learn", that is, methods that leverage data to improve performance on some set of tasks.; 4.4000000000000000, it is seen as a part of artificial intelligence.; 8.0000000000000000, machine learning algorithms build a model based on sample data, known as training data, in order to make predictions or decisions without being explicitly programmed to do so.; 4.0000000000000000, machine learning algorithms are used in a wide variety of applications, such as in medicine, email filtering, speech recognition, and computer vision, where it is difficult or unfeasible to design conventional algorithms to perform the needed tasks.; 2.0000000000000000, A subset of machine learning is closely related to computational statistics, which focuses on making predictions using computers, but not all machine learning is statistical learning.; 4.0000000000000000, The study of mathematical optimization delivers methods, theory and application domains to the field of machine learning.; 2.0000000000000000, Data mining is a related field of study, focusing on exploratory data analysis through unsupervised learning.; 1.0000000000000000, Some implementations of machine learning use data and neural networks in a way that mimics the working of a biological brain.; 2.0000000000000000, In its application across business problems, machine learning is also referred to as predictive analytics.; 2.0000000000000000, Machine learning and data mining often employ the same methods and overlap significantly, but while machine learning focuses on prediction, based on known properties learned from the training data, data mining focuses on the discovery of (previously) unknown properties in the data (this is the analysis step of knowledge discovery in databases); 0.0000000000000000, data mining uses many machine learning methods, but with different goals; on the other hand, machine learning also employs data mining methods as "unsupervised learning" or as a preprocessing step to improve learner accuracy.; 2.0000000000000000, Much of the confusion between these two research communities (which do often have separate conferences and separate journals, KDD/PKDD being a major exception) comes from the basic assumptions they work with: in machine learning, performance is usually evaluated with respect to the ability to reproduce known knowledge, while in knowledge discovery and data mining (KDD) the key task is the discovery of previously unknown knowledge.; 1.0000000000000000, Evaluated with respect to known knowledge, an unsupervised method will usually be outperformed by other supervised methods, while in a typical KDD task, supervised methods cannot be used due to the unavailability of training data.; 1.0000000000000000, Machine learning also has intimate ties to optimization: many learning problems are formulated as minimization of some loss function on a training set of samples.; 0.0000000000000000, Loss functions express the discrepancy between the predictions of the model being trained and the actual problem instances (for example, in classification, one wants to assign a label to instances, and models are trained to correctly predict the pre-assigned labels of a set of examples); 2.0000000000000000 }
```

Vediamo che alcune frasi hanno uno score molto elevato pari a 9.49, il che vuol dire che all'interno del paragrafo vi e' un maggior numero di parole chiave rilevanti per la summarization. Altri paragrafi hanno punteggi molto bassi come 0.22 perche' considerate prive di informazioni utili.

I risultati ottenuti portano alla seguente compressione:

```
machine learning and data mining often employ the same methods and overlap significantly, but while machine learning focuses on prediction, based on known properties learned from the training data, data mining focuses on the discovery of (previously) unknown properties in the data (this is the analysis step of knowledge discovery in databases); data mining uses many machine learning methods, but with different goals; on the other hand, machine learning also employs data mining methods as "unsupervised learning" or as a preprocessing step to improve learner accuracy; Much of the confusion between these two research communities (which do often have separate conferences and separate journals, KDD/PKDD being a major exception) comes from the basic assumptions they work with: in machine learning, performance is usually evaluated with respect to the ability to reproduce known knowledge, while in knowledge discovery and data mining (KDD) the key task is the discovery of previously unknown knowledge; A subset of machine learning is closely related to computational statistics, which focuses on making predictions using computers, but not all machine learning is statistical learning; machine learning algorithms build a model based on sample data, known as training data, in order to make predictions or decisions without being explicitly programmed to do so; Machine learning (ML) is a field of inquiry devoted to understanding and building methods that "learn" - that is, methods that leverage data to improve performance on some set of tasks.
```

Dove sono state riportate le frasi con lo score piu' alto.

Infine e' stata utilizzata una metrica Rouge per poter visualizzare gli effettivi punteggi ottenuti dal task di summarization. Tale metrica consiste in un insieme di metriche per valutare la sintesi automatica dei testi e le traduzioni automatiche. Generalmente funzioni confrontando un riassunto o una traduzione prodotta automaticamente con una serie di summaries (tipicamente prodotte dall'uomo). Le metriche coinvolte in questo processo di valutazione sono:

- Recall
- Precision
- F-1 score

```
rouge_score = [{"rouge-1": {"P": 1.0, "R": 0.3, "F": 0.46153847981677}, "rouge-2": {"P": 0.0075, "R": 0.04442706060606, "F": 0.00571648814353}, "rouge-l": {"P": 1.0, "R": 0.3, "F": 0.46153847981677}}]
```

Vogliamo ricordare che *Rouge-1*, *Rouge-2* e *Rouge-l* misurano l'overlap di unigrammi, bigrammi ed l-grammi presenti nel riassunto e nel documento originale. Come possiamo notare per la metrica Rouge, man mano che aumentiamo il tasso di compressione, subisce una variazione piu' o meno significativa. Poiche' e' la recall che misura il numero di parole sovrapposte in relazione al numero totali di parole presente nel documento originale, questa risulta essere la misura piu' indicativa della metrica.

6 Settima Esercitazione

6.1 Introduzione

La settima esercitazione prevede l'implementazione del *topic modelling*, ovvero creare un modello statistico in grado di determinare gli argomenti o topic di una collezione di documenti. L'esperimento e' stato svolto andando ad estrapolare un corpus da Sketch Engine; l'idea era quella di andare ad estrapolare un corpus che potesse avere una struttura e una composizione organizzata in documenti e paragrafi.

6.2 Sviluppo

Per lo sviluppo sono state implementate le seguenti funzionalita' :

- Solite operazioni di Text Cleaning e preprocessing svolte come nelle esercitazioni del blocco 1.
- Operazione di lettura del corpus: dato il mio file, leggo ogni documento e ogni paragrafo. In output avre' una lista di documenti a cui viene associata una lista di parole pre-processate.
- data la lista di documenti ottenuti dalla funzione precedente, possiamo procedere con lo sviluppo del topic modelling: creo, mediante l'ausilio di *corpora*, un dizionario che contiene il mapping tra le parole e i loro ID interi. A questo dizionario rimuovo tutti i token che non appaiono in meno di 3 documenti e tutti quei token che non appaiono in almeno il 60% dei documenti. Viene poi creata un lista di BoW per ogni termine nel documento. Infine, con *LdaModel* effettuo il training del modello *LDA* (*Latent Dirichlet Allocation*) il cui training e' effettuato a partire da:
 - *corpus_freq* : lista di BoW calcolata precedentemente
 - *num_topics* : numero di topics che verranno estratti dal corpus
 - *id2word*: rappresenta il mio dizionario precedentemente ottenuto ed e' usato per determinare la dimensione del vocabolario e per la stampa del topic

- *passes*: numero di volte che attraverso il corpus in fase di training
- *alpha* : probabilita' a priori per ogni topic
- *chunksize*: numero di documenti da caricare in memoria alla volta.
Il valore di default per chunksize=2000.

6.3 Risultati Ottenuti

Il corpus travel contiene 100 documenti dove, grazie al modello LDA ottenuto, abbiamo ottenuto i seguenti topics:

```
Topic: 0 :
["0.042*class" + 0.018*conditional + 0.018*result + 0.017*third + 0.014*example + 0.012*happy + 0.012*condition + 0.011*situation + 0.010*perfect +
0.009*suppose"]
Topic: 1 :
["0.012*word + 0.010*class + 0.017*man + 0.012*team + 0.011*subject + 0.011*last + 0.010*process + 0.009*question + 0.009*language + 0.009*adjective"]
Topic: 2 :
["0.009*students + 0.009*exam + 0.014*book + 0.012*sb + 0.011*goal + 0.009*lesson + 0.010*speaking + 0.010*language + 0.011*unit + 0.011*read"]
Topic: 3 :
["0.010*learn + 0.011*language + 0.012*lesson + 0.011*word + 0.010*student + 0.009*example + 0.008*grammar + 0.008*situation + 0.008*structure + 0.008*help"]
Topic: 4 :
["0.009*book + 0.011*note + 0.010*book + 0.010*worksheet + 0.010*holiday + 0.014*student + 0.011*test + 0.009*online + 0.008*write + 0.008*topic"]
Topic: 5 :
["0.009*travel + 0.009*article + 0.011*hotel + 0.010*place + 0.011*holiday + 0.011*word + 0.011*trip + 0.010*man + 0.010*stay + 0.009*journey"]
Topic: 6 :
["0.011*content + 0.008*work + 0.008*traffic + 0.008*bus + 0.008*last + 0.007*area + 0.007*vehicle + 0.007*train + 0.007*travel + 0.006*bag + 0.006*place"]
Topic: 7 :
["0.008*rel + 0.017*level + 0.011*video + 0.011*reach + 0.010*student + 0.014*grammar + 0.014*love + 0.014*team + 0.011*lesson + 0.012*class"]
Topic: 8 :
["0.017*travel + 0.014*help + 0.012*example + 0.010*activity + 0.010*music + 0.010*work + 0.009*content + 0.009*phrase + 0.009*post + 0.008*learn"]
```

Notiamo che ogni Topic e' descritto mediante una serie di termini che sono ordinati in base alla loro significativita' all'interno del topic stesso.

Vediamo poi i topics appartenenti ai primi 20 documenti:

```
Doc 0 : [(5, 0.98751235)]
Doc 1 : [(2, 0.13711767), (5, 0.8607147)]
Doc 2 : [(0, 0.5192168), (1, 0.31337816), (3, 0.017647326), (4, 0.13945827)]
Doc 3 : [(1, 0.06805676), (2, 0.39220524), (5, 0.44496346), (7, 0.091738485)]
Doc 4 : [(3, 0.9985373)]
Doc 5 : [(2, 0.01055066), (3, 0.8591707), (9, 0.12990437)]
Doc 6 : [(3, 0.99202764)]
Doc 7 : [(0, 0.22088102), (3, 0.33655864), (5, 0.43885994)]
Doc 8 : [(5, 0.033344958), (9, 0.96604043)]
Doc 9 : [(0, 0.025432905), (5, 0.6630382), (9, 0.30522987)]
Doc 10 : [(0, 0.9990297)]
Doc 11 : [(0, 0.6552332), (4, 0.34283528)]
Doc 12 : [(4, 0.99257296)]
Doc 13 : [(7, 0.99716973)]
Doc 14 : [(2, 0.19703244), (7, 0.7341375), (8, 0.06538537)]
Doc 15 : [(0, 0.07200433), (2, 0.021269456), (3, 0.90590125)]
Doc 16 : [(0, 0.500172), (3, 0.49519813)]
Doc 17 : [(0, 0.26614466), (5, 0.17361982), (9, 0.5588081)]
Doc 18 : [(0, 0.34682715), (3, 0.65098876)]
Doc 19 : [(0, 0.99726015)]
Doc 20 : [(0, 0.30688792), (3, 0.011119908), (9, 0.68103606)]
```