

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 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 indipendente 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 con l'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_FULL.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 3:

- 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 *definitions_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 e antonimi tramite i metodi *total_hyponym_path(word)* (con rispettiva ricerca della distanza dall'iperonimo e la radice che rappresenta il concetto) e *find_antonyms(word)*.
- Infine ho implementato un metodo aggiuntivo, definito *similarity_search(word)* che indica uno score di quanto sono simili le definizioni relativi agli iperonimi rispetto alle definizioni per i concetti analizzati. Per lo svolgimento e' stato utilizzato il modulo *sentence_transformer* che permette di calcolare rappresentazioni vettoriali che riguardano frasi, parafrasi ecc.. Fornisce una serie di modelli pre-addestrati dove la cui frase (o il testo) viene incorporato all'interno dello spazio vettoriale in modo che le frasi simili siano vicini tra loro e che venga calcolata con velocita' la similarita' del coseno. Il modello utilizzato per la rappresentazione vettoriale e' il *all-MiniLM-L6-v2*. Il valore restituito sara' uno per ogni synset del concetto in analisi:

essi corrispondono al valore medio di somiglianza tra le definizioni di tutti gli iperonimi e la definizione della parola in analisi.

2.3 Risultati Ottenuti

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

```
Concept: Paper
Term Synset('paper.n.01') : [['a material made of cellulose pulp derived mainly from wood or rags or certain grasses', 15]]
Term Synset('composition.n.00') : [['an essay (especially one written as an assignment)', 8]]
Term Synset('newspaper.n.01') : [['a daily or weekly publication on folded sheets; contains news and articles and advertisements', 14]]
Term Synset('paper.n.04') : [['a medium for written communication', 5]]
Term Synset('paper.n.05') : [['a scholarly article describing the results of observations or stating hypotheses', 11]]
Term Synset('newspaper.n.02') : [['a business firm that publishes newspapers', 6]]
Term Synset('newspaper.n.03') : [['the physical object that is the product of a newspaper publisher', 11]]
Term Synset('paper.v.01') : [['cover with paper', 3]]
Term Synset('wallpaper.v.01') : [['cover with wallpaper', 3]]
```

Come possiamo vedere dalla figura sopra, per ogni concetto che analizziamo, ne riportiamo la rispettiva lunghezza della definizione presa da WordNet:

Noun

- [S: \(n\)](#) **paper** (a material made of cellulose pulp derived mainly from wood or rags or certain grasses)
- [S: \(n\)](#) **composition, paper, report, theme** (an essay (especially one written as an assignment)) *"he got an A on his composition"*
- [S: \(n\)](#) **newspaper, paper** (a daily or weekly publication on folded sheets; contains news and articles and advertisements) *"he read his newspaper at breakfast"*
- [S: \(n\)](#) **paper** (a medium for written communication) *"the notion of an office running without paper is absurd"*
- [S: \(n\)](#) **paper** (a scholarly article describing the results of observations or stating hypotheses) *"he has written many scientific papers"*
- [S: \(n\)](#) **newspaper, paper, newspaper publisher** (a business firm that publishes newspapers) *"Murdoch owns many newspapers"*
- [S: \(n\)](#) **newspaper, paper** (the physical object that is the product of a newspaper publisher) *"when it began to rain he covered his head with a newspaper"*

Verb

- [S: \(v\)](#) **paper** (cover with paper) *"paper the box"*
- [S: \(v\)](#) **wallpaper, paper** (cover with wallpaper)

Per quanto riguarda il secondo studio, vogliamo mostrare i risultati che riguardano la lunghezza della definizione per ciascun iperonimo del concetto analizzato e anche la rispettiva distanza tra la radice e il termine :

```
Concept: courage
hypernym distance from the root of the concept Synset('courage.n.01'):
((Synset('character.n.03'), 2), (Synset('courage.n.01'), 0), (Synset('spirit.n.03'), 1), (Synset('trait.n.01'), 3), (Synset('entity.n.01'), 6), (Synset('abstraction.n.06'), 11), (Synset('attribute.n.02'), 9), (Synset('trait.n.01'), 7), (Synset('character.n.03'), 2))
hypernym length definition :
((Synset('entity.n.01'), 17), (Synset('abstraction.n.06'), 11), (Synset('attribute.n.02'), 9), (Synset('trait.n.01'), 7), (Synset('character.n.03'), 2))
```

Stesso ragionamento e' stato fatto per la ricerca degli antonimi, anche se ne esiste solo uno per il concetto *Courage*:

```

-----
Concept:  Courage
Antonym  : [(Lemma('cowardice.n.01.cowardice'), 1)]

-----

Concept:  Paper
Antonym  : []

-----

Concept:  Apprehension
Antonym  : []

-----

Concept:  Sharpener
Antonym  : []

```

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

```

Noun
• 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
  ◦ antonym
    • W: (n) cowardice [Opposed to: courage] (the trait of lacking courage)
  ◦ derivationally related form

```

Infine vogliamo mostrare i risultati ottenuti dall'ultimo studio:

```

Concept: Paper
concetto Synset('paper.n.01') :
Average Similarity for hypernyms : 0.22500497102737427

concetto Synset('composition.n.08') :
Average Similarity for hypernyms : 0.4687023162841797

concetto Synset('newspaper.n.01') :
Average Similarity for hypernyms : 0.5600683093070984

concetto Synset('paper.n.04') :
Average Similarity for hypernyms : 0.49245941638946533

concetto Synset('paper.n.05') :
Average Similarity for hypernyms : 0.23190630972385406

concetto Synset('newspaper.n.02') :
Average Similarity for hypernyms : 0.7873624563217163

concetto Synset('newspaper.n.03') :
Average Similarity for hypernyms : 0.31530842185020447

concetto Synset('paper.v.01') :
Average Similarity for hypernyms : 0.5753957033157349

concetto Synset('wallpaper.v.01') :
Average Similarity for hypernyms : 0.40710264444351196

```

Dove vediamo che solo alcune definizioni degli iperonimi sono mediamente vicine alla definizione del concetto *Paper* analizzato in questione.

3 Quarta Esercitazione

3.1 Introduzione

La quarta esercitazione prevede lo studio di alcuni verbi dal punto di vista della Teoria di Hank. Secondo Hank il verbo e' 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 e' formato da un certo numero di valori che lo riempiono, detti *filler*. Ogni filler puo' avere associati dei tipi semantici che rappresentano delle generalizzazioni strutturate come una gerarchia.

3.2 Sviluppo

Per lo svolgimento e' stato utilizzato il seguente verbo :

- *Read, Reads, Readen*

Per recuperare le n frasi in cui il verbo viene utilizzato, e' stato utilizzato il corpus brown: vengono recuperate una serie di frasi dal corpus che contengono o il verbo presente *read* e *reads*, oppure la sua forma al present perfect *readen*. Poiche' si tratta di recuperare una serie di frasi dal corpus, lavoreremo con

una serie di contesti diversi a seconda del ruolo che viene svolto dal verbo in ciascuna frase. Per fare cio', applico l'algoritmo di lesk che rappresenta uno dei cardini della Word Sense Disambiguation. L'algoritmo di Lesk si basa sul fatto che, all'interno di una regione testuale, un insieme di parole condividono un significato simile. Verra' restituito il *best_sense* per una determinata parola a partire da un insieme di contesti presi da esempi per quella parola.

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, e' possibile estrarre il soggetto e il complemento oggetto relativi al verbo di una frase che vengono poi mappati con i relativi super sensi di WordNet. 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 piu' frequenti.

3.3 Risultati Ottenuti

Vediamo che prima di tutto possiamo mostrare una parte del soggetto e del complemento oggetto estratto da ciascuna frase:

```
[ 'You', 'read', 'various', 'guesses', 'on', 'how', 'many', 'Americans', 'are', 'camping', '.' ]
sub : You      obj: guesses
[ 'Again', ',', 'the', 'oil', 'men', 'must', 'read', 'the', 'motors', 'at', 'such', 'intervals', 'as', 'he', 'finds', 'best', '.' ]
sub : men      obj: actors
[ 'She', 'read', 'Hailand's', 'back', 'Ages', '...', '...', 'which', 'it', 'enjoyed', 'very', 'much', '...', '...', '...' ]
sub : She      obj: Ages
```

La cui coppia di cluster semantici e' pari a :

```
Coppie Semantiche with count:
[('noun.substance', 'noun.artifact'), 1]
[('noun.person', 'adj.all'), 1]
[('noun.substance', 'noun.communication'), 4]
[('noun.substance', 'adv.all'), 1]
[('noun.person', 'noun.cognition'), 2]
[('noun.person', 'verb.social'), 1]
[('noun.location', 'noun.cognition'), 1]
[('adj.all', 'noun.group'), 1]
[('noun.person', 'noun.group'), 1]
[('noun.cognition', 'noun.person'), 1]
[('noun.quantity', 'noun.communication'), 2]
[('noun.substance', 'noun.attribute'), 1]
[('noun.person', 'noun.quantity'), 2]
[('noun.person', 'noun.time'), 1]
[('noun.person', 'noun.person'), 1]
[('noun.person', 'noun.communication'), 5]
[('noun.substance', 'noun.location'), 1]
[('noun.person', 'noun.act'), 1]
[('noun.substance', 'noun.person'), 1]
[('noun.quantity', 'noun.cognition'), 1]
[('noun.person', 'adv.all'), 1]
```

dove possiamo notare che per il verbo *read*, abbiamo che la coppia di cluster semantici piu' frequente e' quella relativa alla coppia *person:communication*, con una frequenza pari a 5.

4 Quinta Esercitazione

4.1 Introduzione

La quinta esercitazione prevede la sperimentazione del *content-to-form*, cioe' cercare di risalire al synset di un concetto indirizzando la ricerca in WordNet attraverso i genus. Alla base, vi e' il principio *Genus-Differentia definition*, secondo il quale un concetto puo' 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 piu' 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 piu' importanti da estrarre per un concetto. La funzione mi restituira' una lista ordinata dei 5 termini piu' importanti (ovvero quelli piu frequenti).
- Per ogni genus ottengo un possibile candidato; ottengo una lista di iponimi di tutti i synset relativi al termine genus. Verra' restituito quel synset che si avvicina di piu' 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('holoid.n.01'))]
-----
Concept: Paper
Genus list (with frequency):
[('material', 23), ('write', 10), ('cellulose', 7), ('wood', 6), ('tree', 5)]

Candidates:
[('material', Synset('composite_material.n.01')), ('write', Synset('handwrite.v.01')), ('cellulose', Synset('pulp.n.03')), ('wood', Synset('balsa.n.01')), ('tree', Synset('poon.n.02'))]
-----
Concept: Apprehension
Genus list (with frequency):
[('fear', 10), ('anxiety', 10), ('feeling', 5), ('happen', 5), ('feel', 4)]

Candidates:
[('fear', Synset('apprehension.n.01')), ('anxiety', Synset('paria.n.02')), ('feeling', Synset('glow.v.05')), ('happen', Synset('concur.v.02')), ('feel', Synset('glow.v.05'))]
-----
Concept: Sharpener
Genus list (with frequency):
[('pencil', 25), ('sharpen', 17), ('tool', 10), ('object', 11), ('allow', 4)]

Candidates:
[('pencil', Synset('lead_pencil.n.01')), ('sharpen', Synset('edge.v.04')), ('tool', Synset('drill.n.01')), ('object', Synset('commemorative.n.01')), ('allow', Synset('pass.v.17'))]

```

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 puo' essere dovuto da un legame semantico debole che c'e' tra il concetto preso in analisi e l'intorno ottenuto attraverso i genus e i lori iperonimi. Inoltre, il documento *defs.csv* e' un file che e' stato rappresentato da noi studenti, quindi alcune definizioni possono essere lontane da cio' 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 rappresenta 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 stopwords, punteggiatura, oltre grazie al supporto fornito dalla libreria spacy.
- Successivamente, posso passare all'operazione di lettura del mio documento. Per fare cio' 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 piu' frequenti all'interno del mio file: la frequenza delle parole mi sara' 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 e' 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 visualizzarne 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 :

```
Sentences : (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.: 4.055555555555555, It is seen as a part of artificial intelligence.: 0.2222222222222222, 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.: 4.277777777777777, 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 develop conventional algorithms to perform the needed tasks.: 2.344444444444444, 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.: 4.888888888888889, the study of mathematical optimization delivers methods, theory and application domains to the field of machine learning.: 2.666666666666667, data mining is a related field of study, focusing on exploratory data analysis through unsupervised learning.: 3.444444444444444, Some implementations of machine learning use data and neural networks in a way that mimics the working of a biological brain.: 2.666666666666667, In its application across business problems, machine learning is also referred to as predictive analytics.: 2.166666666666667, 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): 9.499999999999999, 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.: 7.999999999999999, Much of the confusion between these two research communities (which do often have separate conferences and separate journals, ECDM 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.: 5.833333333333333, evaluated with respect to known knowledge, an uninformed (unsupervised) method will easily be outperformed by other supervised methods, while in a typical KDD task, supervised methods cannot be used due to the unavailability of training data.: 3.444444444444444, Machine learning also has intimate ties to optimization: many learning problems are formulated as minimization of some loss function on a training set of examples.: 3.833333333333326, 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): 2.0000000000000004)
```

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, ECDM 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": {"r": 1.0, "p": 0.3, "f": 0.4615384579881657}, "rouge-2": {"r": 0.6875, "p": 0.034482758620689655, "f": 0.06567164088144353}, "rouge-l": {"r": 1.0, "p": 0.3, "f": 0.4615384579881657}}]
```

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 una 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**"clause" + 0.018**"conditionals" + 0.018**"result" + 0.017**"third" + 0.014**"example" + 0.012**"happen" + 0.012**"condition" + 0.011**"situation" + 0.010**"perfect" +
0.010**"student"]
Topic 1 :
["0.022**"word" + 0.018**"clause" + 0.017**"noun" + 0.013**"tense" + 0.011**"subject" + 0.011**"learn" + 0.010**"pronoun" + 0.010**"question" + 0.009**"language" + 0.009**"adjective"]
Topic 2 :
["0.050**"students" + 0.029**"exam" + 0.024**"book" + 0.022**"is" + 0.021**"speak" + 0.020**"lesson" + 0.019**"speaking" + 0.018**"language" + 0.017**"unit" + 0.017**"read"]
Topic 3 :
["0.011**"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 :
["0.029**"book" + 0.021**"hotel" + 0.020**"back" + 0.018**"worksheet" + 0.016**"holiday" + 0.014**"student" + 0.012**"test" + 0.009**"online" + 0.008**"write" + 0.008**"topic"]
Topic 5 :
["0.028**"travel" + 0.020**"article" + 0.017**"hotel" + 0.016**"place" + 0.015**"holiday" + 0.013**"word" + 0.011**"trip" + 0.010**"noun" + 0.010**"stay" + 0.009**"journey"]
Topic 6 :
["0.011**"content" + 0.008**"work" + 0.008**"traffic" + 0.008**"bus" + 0.007**"turn" + 0.007**"vehicle" + 0.007**"train" + 0.007**"travel" + 0.006**"buy" + 0.006**"place"]
Topic 7 :
["0.020**"oil" + 0.022**"level" + 0.021**"video" + 0.021**"teach" + 0.018**"student" + 0.014**"grammar" + 0.014**"love" + 0.014**"tense" + 0.013**"lesson" + 0.012**"clause"]
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**"phrasal" + 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)]
```