



The banner features a row of six icons: a globe, a book, a handshake, a money bag with a Euro symbol, a scale of justice, and a bicycle. Below the icons, the text reads 'AIUCD 2021' in large black letters, followed by 'DH per la società: e-guaglianza, partecipazione, diritti e valori nell'era digitale' and '10° congresso annuale PISA 19-22 gennaio'. On the right side, a list of topics is displayed in colored text: 'DIGITAL PUBLIC HUMANITIES' (red), 'OPEN CULTURE' (orange), 'RETI SOCIALI' (yellow), 'TECH ECONOMY' (green), 'E-PARTICIPATION' (blue), and 'TECNOLOGIE ASSISTIVE' (purple). The background includes binary code and a classical building facade.

AIUCD 2021

DH per la società: e-guaglianza, partecipazione, diritti e valori nell'era digitale

10° congresso annuale **PISA** 19-22 gennaio

DIGITAL PUBLIC HUMANITIES
OPEN CULTURE
RETI SOCIALI
TECH ECONOMY
E-PARTICIPATION
TECNOLOGIE ASSISTIVE

Versione PROVVISORIA del contributo presentato al Convegno Annuale

DISCLAIMER

Questa versione dell'abstract non è da considerarsi definitiva e viene pubblicata esclusivamente per facilitare la partecipazione del pubblico al convegno AIUCD 2021

Il Book of Abstract contenente le versioni definitive e dotato di ISBN sarà disponibile liberamente a partire dal 19 gennaio sul sito del convegno sotto licenza creative commons.

Analisi del sentiment delle Confessioni di Sant'Agostino

Angelica Lo Duca¹, Andrea Marchetti¹

¹ Institute of Informatics and Telematics, National Research Council
via G. Moruzzi 1, 56124 Pisa (Italy)
[name].[surname]@iit.cnr.it

ABSTRACT

Le Confessioni di Sant'Agostino costituiscono uno dei capolavori della letteratura cristiana di ogni tempo. In quest'opera, costituita da 13 libri, il celebre scrittore descrive il proprio percorso di conversione a Dio. Il presente articolo applica tecniche di sentiment analysis a ciascun libro, attraverso due tecniche, una basata sull'Apprendimento Supervisionato e l'altra sull'Apprendimento Non Supervisionato, al fine di capire lo stato d'animo dell'autore man mano che il percorso di conversione arriva a compimento. Come risultato, si può evincere che mentre nei primi libri i sentimenti positivi sono comparabili a quelli negativi, negli ultimi libri, i sentimenti negativi tendono totalmente a scomparire.

PAROLE CHIAVE

Sentiment Analysis, Digital Humanities, Sant'Agostino, Confessioni.

1. INTRODUZIONE

Le Confessioni di Sant'Agostino sono uno dei maggiori capolavori della letteratura cristiana. Opera scritta tra il 397 e il 400 d.C., le Confessioni riflettono la vita dell'autore. Esse si suddividono in 13 libri, nei quali l'autore racconta il proprio cammino di conversione a Dio, descrivendo gli eventi che lo hanno portato a diventare cristiano, a partire dall'infanzia. Mentre i primi nove libri sono totalmente autobiografici, gli ultimi quattro riguardano diverse concezioni filosofiche, come l'origine dell'uomo e l'essenza del tempo. L'evento della conversione avviene propriamente nel nono libro.

L'obiettivo di questo lavoro consiste nell'analisi del sentiment nell'ambito dei 13 libri che costituiscono le Confessioni di Sant'Agostino, al fine di stabilire se esista o meno una relazione con il percorso della conversione¹. Nello specifico, sono descritti due esperimenti di analisi del sentiment, uno basato su tecniche di Apprendimento Supervisionato e l'altro su tecniche di Apprendimento Non Supervisionato. Sebbene i risultati delle due tecniche siano leggermente diversi, si conferma il trend che negli ultimi due libri, prevale il sentiment positivo. Ciò è in linea con il processo di conversione descritto nel corso dei libri.

2. STATO DELL'ARTE

Il problema della Sentiment Analysis di un testo può essere affrontato sia utilizzando tecniche di Apprendimento Supervisionato sia tecniche di Apprendimento Non Supervisionato [1]. In generale le prime hanno prestazioni migliori rispetto alle seconde [2]. L'Apprendimento Supervisionato (AS) è una tecnica di Machine Learning (ML) che si basa sull'annotazione iniziale di un dataset di riferimento, al fine di identificare un insieme finito di classi di uscita. Nell'ambito della Sentiment Analysis, le tecniche di AS maggiormente utilizzate sono basate su algoritmi classici di ML, come Support Vector Machines, Naive Bayes [3, 4, 5] e Deep Convolutional Neural Networks [6, 7].

L'Apprendimento Non Supervisionato (ANS) è una tecnica di Machine Learning che cerca di identificare le caratteristiche comuni dei dati forniti in ingresso, al fine di classificare i dati futuri forniti in ingresso. Non è prevista alcuna annotazione iniziale dei dati, in quanto le classi di uscita sono calcolate automaticamente dall'algoritmo. Nell'ambito della Sentiment Analysis, le tecniche di ANS si basano su k-means [8] e LDA [9].

In letteratura esistono altre tecniche per effettuare la Sentiment Analysis di un testo, basate sull'uso di lessici standard, come ad esempio come SentiWordNet [10], AFINN [11], o lessici emotivi, come ad esempio NRC/Emolex [12] e Vader [13], senza contare senza contare i dizionari creati tramite word embeddings [14].

3. METODOLOGIA

Al fine di definire il sentiment di ogni libro, è stata considerata la versione in inglese del testo, fornita dal progetto Gutenberg². Ogni libro è stato analizzato separatamente e diviso in frasi. La Tabella 1 mostra alcune statistiche per ogni

¹ Il software implementato per eseguire le analisi può essere scaricato dal seguente repository Github: <https://github.com/alod83/papers/tree/master/aiucd2021>

² La versione in inglese del testo è disponibile al seguente link: <https://www.gutenberg.org/files/3296/3296-h/3296-h.htm>

libro, mentre la Fig. 1 mostra la word cloud relativa alle parole maggiormente utilizzate nell'intera opera, da cui sono state eliminate le stopwords.

Libro	1	2	3	4	5	6	7	8	9	10	11	12	13
#frasi	186	98	119	182	153	178	185	208	207	455	258	176	303
#caratteri	37735	20500	30011	37217	36606	42798	45414	43764	47943	86515	49733	55939	68395
#parole	7055	3781	5466	6975	6695	7838	8409	8123	8819	16315	9501	10205	12598

Tabella 1. Statistiche per ogni libro.

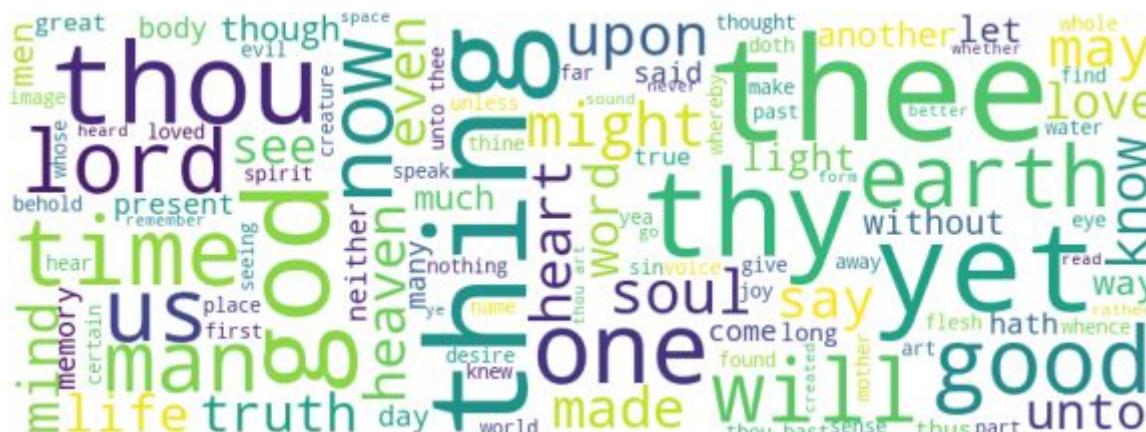


Fig. 1 Word Cloud relativa alle parole maggiormente utilizzate nell'ambito dell'intera opera delle Confessioni.

Per ogni frase sono stati calcolati il sentiment positivo, negativo e neutro (come specificato in seguito). A partire dai valori dei sentiment ottenuti, sono stati calcolati gli indici di negatività, positività e neutralità di ciascun libro i , come segue:

$$neg_i = \frac{\text{numero di frasi nel libro } i \text{ con un sentiment negativo}}{\text{numero totale di frasi nel libro } i}$$

$$pos_i = \frac{\text{numero di frasi nel libro } i \text{ con un sentiment positivo}}{\text{numero totale di frasi nel libro } i}$$

$$neu_i = \frac{\text{numero di frasi nel libro } i \text{ con un sentiment neutrale}}{\text{numero totale di frasi nel libro } i}$$

Il sentiment relativo ad ogni libro è stato calcolato utilizzando due tecniche diverse, una basata sull'Apprendimento Supervisionato e una sull'Apprendimento Non Supervisionato.

4. ESPERIMENTI

Nell'ambito dell'ANS, è stato utilizzato il lessico AFINN, nella versione disponibile per python³. Questo software, che supporta solo l'inglese e il danese, prende in ingresso una frase e restituisce in uscita un valore (score) positivo, negativo o neutro. Lo score ottenuto corrisponde al sentiment della frase. Quindi una frase con lo score maggiore di zero corrisponde ad una frase con sentiment positivo. Similmente, una frase con lo score minore di zero corrisponde ad una frase con sentiment negativo. Infine, una frase con lo score uguale a zero corrisponde ad una frase con sentiment nullo. A partire dallo score di ogni frase, sono stati calcolati gli indici di negatività, positività e neutralità. La Fig. 2 mostra il risultato del processo di analisi.

³ <https://pypi.org/project/afinn/>

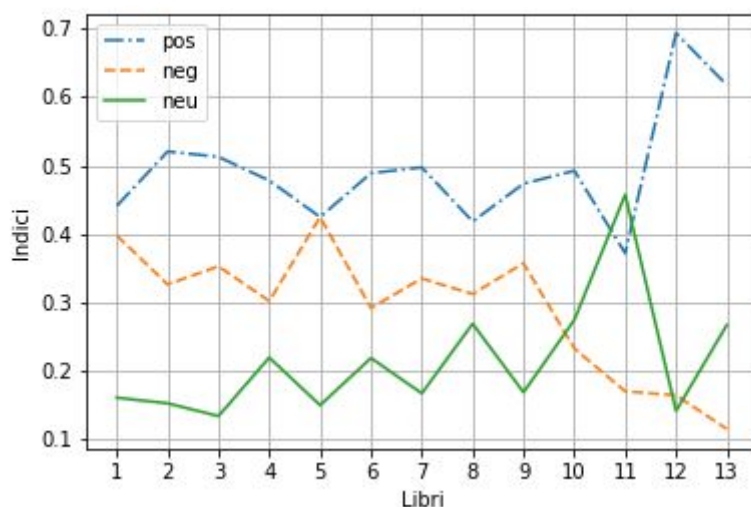


Fig. 2 Indici di positività, negatività e neutralità nei vari libri, utilizzando il lessico AFINN.

Nell'ambito dell'AS, è stato utilizzato il seguente procedimento: ogni frase del primo libro è stata annotata, attribuendo un sentiment positivo, nullo, o negativo. In totale, sono state ottenute 82 frasi con sentiment positivo, 74 con sentiment negativo e 30 con sentiment nullo. Del totale, sono state considerate solo le frasi con sentiment positivo e negativo, in modo da costruire il dataset di riferimento (156 campioni), utilizzato per allenare l'algoritmo di AS. La Tabella 2 mostra alcuni campioni del dataset di riferimento: la frase costituisce l'input e il sentiment l'output dell'algoritmo.

Frase	Sentiment
Great art Thou, O Lord, and greatly to be praised; great is Thy power, and Thy wisdom infinite.	1
And Thee would man praise; man, but a particle of Thy creation; man, that bears about him his mortality, the witness of his sin, the witness that Thou resistest the proud: yet would man praise Thee; he, but a particle of Thy creation.	1
For if I go down into hell, Thou art there.	0
I could not be then, O my God, could not be at all, wert Thou not in me; or, rather, unless I were in Thee, of whom are all things, by whom are all things, in whom are all things? Even so, Lord, even so.	1

Tabella 2. Alcuni campioni del dataset usato per allenare l'algoritmo di AS.

Il dataset di riferimento è stato diviso in due parti, training set, contenente il 75% dei campioni e il test set, contenente il rimanente 25%. Ciascuna frase è stata divisa in token al fine di costruire la matrice dei token. Sono stati considerati solo i token con frequenza minima pari a 5 per cui il numero dei nomi contenuti nella matrice dei token è pari a 189. Sono stati condotti due esperimenti, uno senza gli n-grammi, l'altro che considera n-grammi di dimensione 2. Come algoritmo di AS è stato considerato il Linear Support Vector Machines. La Fig. 3 mostra i risultati per l'AS, senza e con l'uso di n-grammi di dimensione 2. Nei due casi i risultati sono abbastanza simili, con una differenza nel valore dello Score AUC, pari a 0.69 nel primo caso e 0.82 nel secondo.

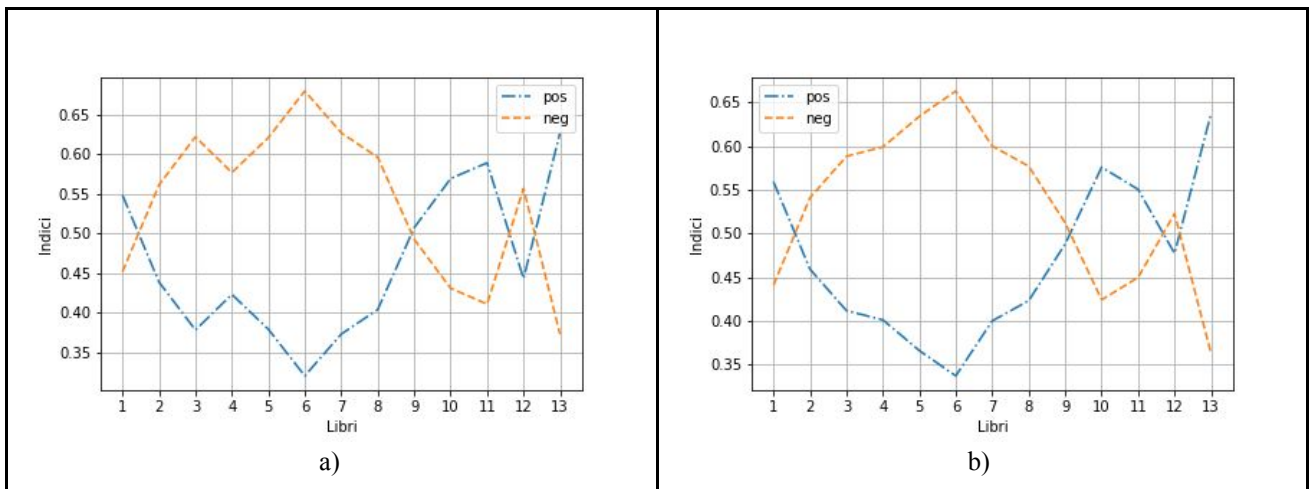


Fig. 3 Indici di positività e negatività nei vari libri, utilizzando l'algoritmo Linear Support Vector Machines, senza (a) e con l'uso di n-grammi (b).

5. DISCUSSIONE E SVILUPPI FUTURI

Ad una prima analisi, è interessante notare come i due esperimenti diano risultati abbastanza diversi. Infatti, mentre nell'esperimento basato su ANS il sentiment positivo è sempre maggiore di quello negativo, nell'esperimento basato su AS, ciò è vero solo nei capitoli 9, 10, 11 e 13. In realtà, entrambi gli esperimenti mostrano come nei primi libri dell'opera il sentiment negativo sia molto alto, mentre negli ultimi libri vada scemando, per lasciare il posto ad un sentiment positivo molto alto. In conclusione, si può dire che entrambi gli esperimenti riflettono a grandi linee il processo di conversione dell'autore, che passa da uno stato d'animo negativo (all'inizio del libro) ad uno stato d'animo positivo (alla fine del libro).

Il presente lavoro costituisce un punto di partenza nell'analisi di questa splendida opera letteraria. Come sviluppo futuro, si potrebbe pensare di analizzare il collegamento tra il sentiment descritto in ogni libro e il contenuto del libro stesso, con un focus particolare alle fasi della vita dell'autore.

BIBLIOGRAFIA

- [1] Kim, E., & Klinger, R. (2018). A survey on sentiment and emotion analysis for computational literary studies. arXiv preprint arXiv:1808.03137.
- [2] Pang B, Lee L, Vaithyanathan S (2002) Thumbs up?: sentiment classification using machine learning techniques. In: Proceedings of the ACL-02 conference on Empirical methods in natural language processing- Volume 10, Association for Computational Linguistics, pp 79–86.
- [3] Al-Hadhrani S, Al-Fassam N, Benhidour H (2019) Sentiment analysis of english tweets: A comparative study of supervised and unsupervised approaches. In: 2019 2nd International Conference on Computer Applications & Information Security (ICCAIS), IEEE, pp 1–5.
- [4] Durant KT, Smith MD (2006) Mining sentiment classification from political web logs. In: Proceedings of Workshop on Web Mining and Web Usage Analysis of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (WebKDD- 2006), Philadelphia, PA.
- [5] Pang B, Lee L, Vaithyanathan S (2002) Thumbs up?: sentiment classification using machine learning techniques. In: Proceedings of the ACL-02 conference on Empirical methods in natural language processing- Volume 10, Association for Computational Linguistics, pp 79–86.
- [6] Jianqiang Z, Xiaolin G, Xuejun Z (2018) Deep convolution neural networks for twitter sentiment analysis. IEEE Access 6:23253–23260.
- [7] Severyn A, Moschitti A (2015) Unitn: Training deep convolutional neural network for twitter sentiment classification. In: Proceedings of the 9th international workshop on semantic evaluation (SemEval 2015), pp 464–469.
- [8] Li N, Wu DD (2010) Using text mining and sentiment analysis for online forums hotspot detection and forecast. Decision support systems 48(2):354–368.
- [9] Li F, Huang M, Zhu X (2010) Sentiment analysis with global topics and local dependency. In: Twenty-Fourth AAAI Conference on Artificial Intelligence.
- [10] Baccianella S, Esuli A, Sebastiani F (2010) Sentiword net 3.0: an enhanced lexical resource for sentiment analysis and opinion mining. In: Lrec, vol 10, pp 2200– 2204.
- [11] Nielsen FÅ (2011) Afinn. Richard Petersens Plads, Building 321.
- [12] Mohammad, S. M., & Turney, P. D. (2013). Crowdsourcing a Word-Emotion Association Lexicon. Computational Intelligence, 29, 436–465.
- [13] Gilbert, C. H. E., & Hutto, E. (2014, June). Vader: A parsimonious rule-based model for sentiment analysis of social media text. In Eighth International Conference on Weblogs and Social Media (ICWSM-14). Available at (20/04/16) [http://comp. social. gatech. edu/papers/icwsm14.vader.hutto.pdf](http://comp.social.gatech.edu/papers/icwsm14.vader.hutto.pdf) (Vol. 81, p. 82).
- [14] Giatsoglou, M., Vozalis, M. G., Diamantaras, K., Vakali, A., Sarigiannidis, G., & Chatzivasvas, K. C. (2017). Sentiment analysis leveraging emotions and word embeddings. Expert Systems with Applications, 69, 214-224.