

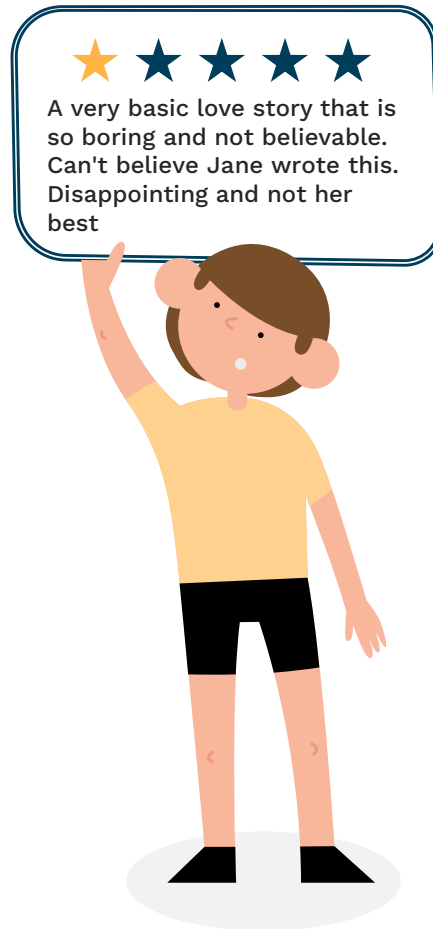
Appello Febbraio 2022

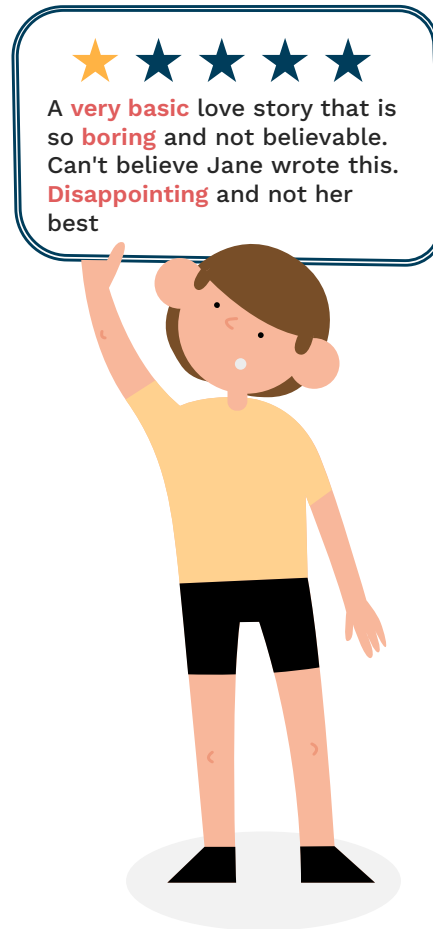
TEXT MINING PROJECT

Sentiment analysis and topic
modelling of Amazon e-books
reviews

Alfredo Galli & Riccardo Rubini









Descrizione **Dataset**

Il dataset utilizzato è stato ottenuto unendo i due file JSON, uno per le recensioni uno per i metadati, provenienti dalla repository *Amazon Review Data (2018)* di *Jianmo Ni*



486 mila il numero totale di recensioni



331 caratteri la lunghezza mediana



4,34 stelle il voto medio



40 categorie relative ai libri

well written first book read book really enjoyed highly recommend even honest review

book serie relationship forward next feeling reading book loved book husband book read enjoyed reading

mind either problem put book often talk child still son add plan use great read good story instead always sure learn hand tell copy book nothing

great story especially think name perfect girl book great children fact rather said sister everything exchange honest wait next although

book one book author looking life realize call best friend parent someone love serie needed world

main character characters well almost work novella daughter moment forward reading guy experience idea now

book good read next really good brother help good book used made book love

high school need enjoyed story Highly recommended secret yet novel series book

really enjoyed highly recommend even though other situation goe good read making book book believe must read story line thinking mother

love story next one thought least become sort may book will family live

great book thought become sort may book will family live

short story next story recommend book quick read beautiful

hero read one will definitely one favorite character development really liked kind love book end book point

maybe wait read look forward first time heroine find well done writing style second book twists turn

Pre-processing workflow



Rimozione
NA



lowercasing



Filtraggio
stopwords



Rimozione
punteggiatura



Rimozione
white space



Stemming

Testo non processato

“I really enjoyed the book. Had the normal back against the wall moments. It even had laugh out loud moments”



Testo processato

“realli enjoy book normal back wall moment even laugh loud moment”

SENTIMENT ANALYSIS



Creazione Variabile Target

La variabile target è stata creata partendo dal rating espresso in stelle nel seguente modo:

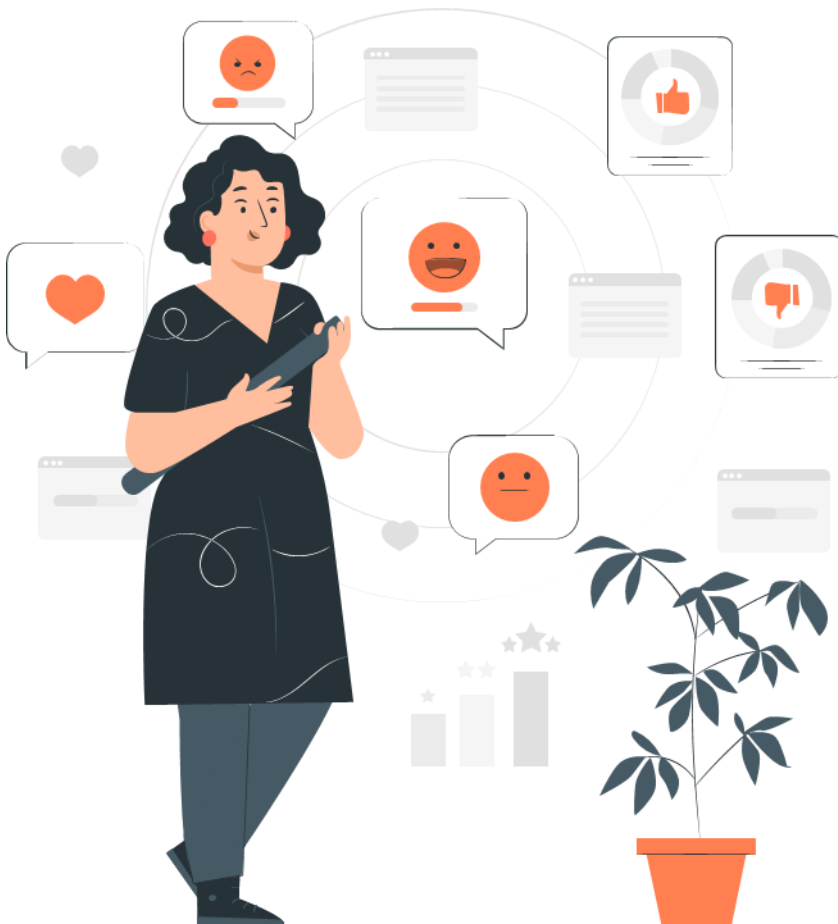
- Si sono **eliminate le recensioni con rating pari a 3** essendo ambigue
- Le recensioni con **valutazione superiore a 4** sono state etichettate come **positive**
- Le recensioni con **valutazione inferiore a 2** sono state etichettate come **negative**

92.1%

La percentuale di recensioni classificate come **positive**

7.9%

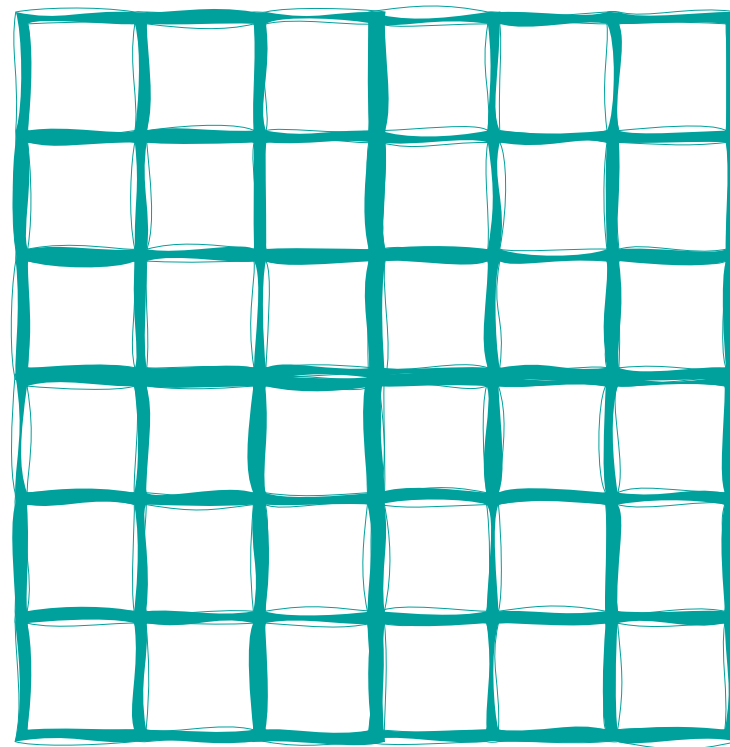
La percentuale di recensioni classificate come **negative**



Recensioni Raw



**Rappresentazione
TF-IDF con bi-grammi**



Feature Selection

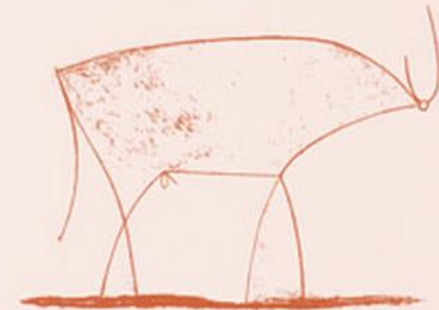
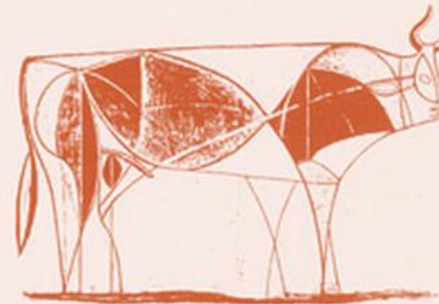
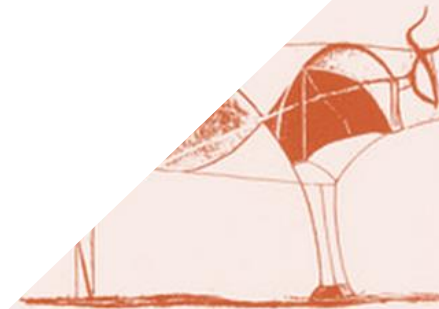
Il processo di feature selection consente di estrarre dal dataset il **set di variabili più significative ed essenziali** per la task di classificazione, in modo da alleggerire lo sforzo computazionale.

Questo passaggio è stato effettuato tramite l'approccio di eliminazione delle feature **NEAR ZERO VARIANCE**, così da eliminare quelle variabili sulla carta dallo scarso valore informativo per il classificatore

**1.5 mln
feature**



**4.8 mila
feature**



Logistic Lasso

Come classificatore la scelta è ricaduta su un **modello logistico penalizzato**, nel particolare con metodo Lasso. Si è fatta questa scelta così da:

- Ottenere **risultati interpretabili**
- Eseguire una **model selection** già direttamente nella fase di training del modello. In totale le feature con coefficiente diverso da zero (i.e. selezionate) sono state 3268

Inoltre, dato la situazione unbalanced nella variabile target, si sono utilizzati pesi in modo da dare più rilevanza alle classificazioni negative corrette, in modo da bilanciare il classificatore.

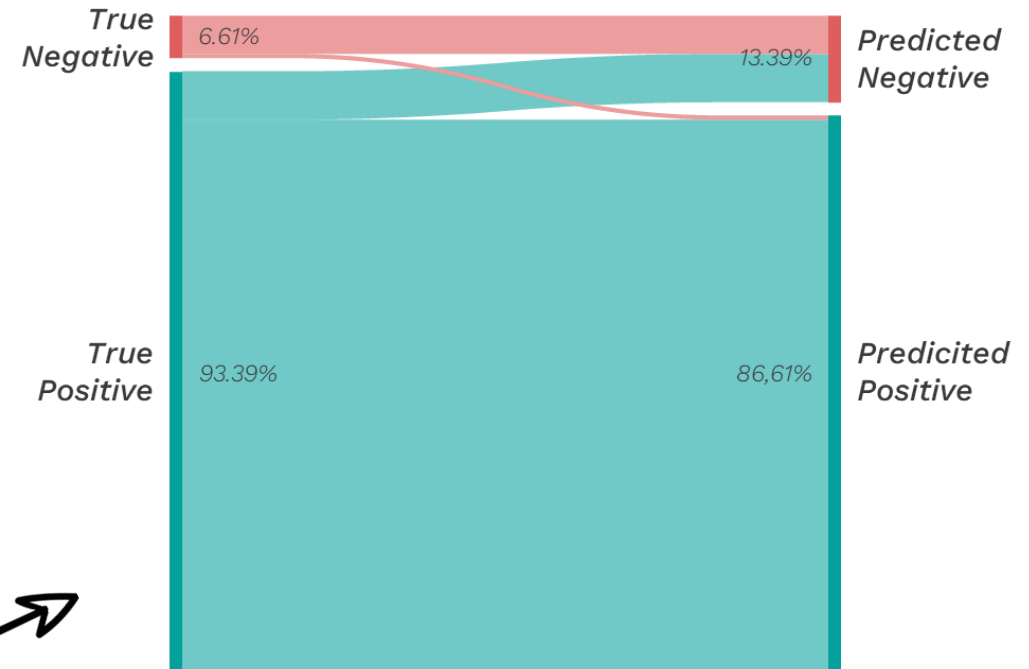


Performance

Test Set

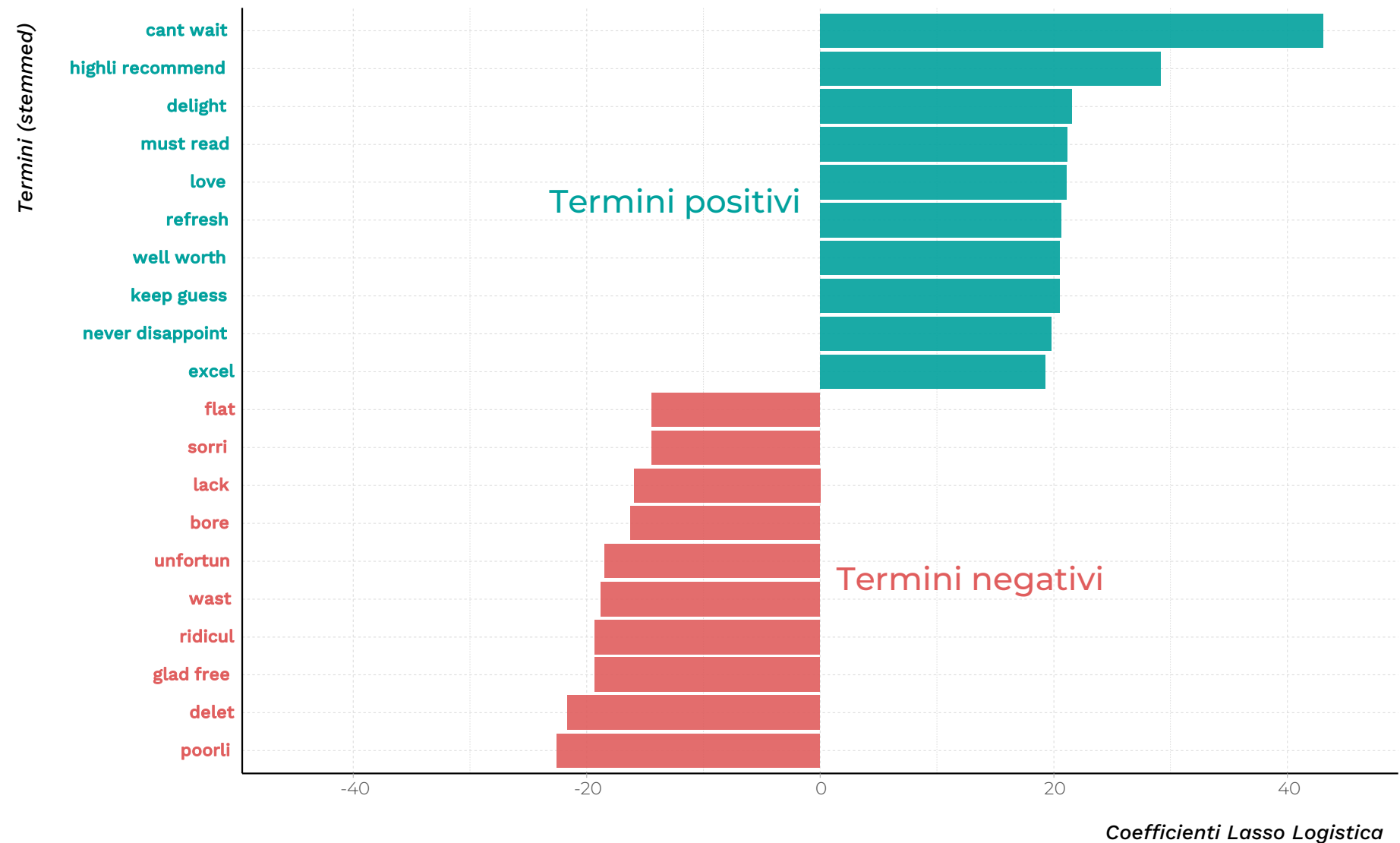
		Predicted class	
		Positive Reviews	Negative Reviews
True class	Positive Reviews	85.95%	7.44%
	Negative Reviews	0.66%	5.95%

Calcolato su 110k recensioni di test



- **ACCURACY:** 92%
- **RECALL POSITIVE:** 92%
- **RECALL NEGATIVE:** 90%

Termini più discriminanti



TOPIC MODELING





Topic Modeling

Con l'applicazione del topic modeling il nostro obiettivo era quello di **estrarre** dalla collezione di recensioni gli **argomenti principali latenti** discussi al loro interno. Per far ciò sono stati utilizzati due approcci:

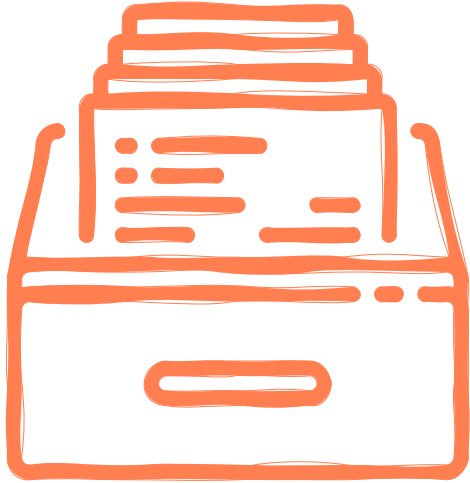
1. **NMF**: approccio puramente **algebrico** di riduzione della dimensionalità della matrice documenti-termini.
2. **LDA**: approccio di tipo **probabilistico**

Per evitare di includere parole troppo rare o troppo popolari, le quali avrebbero potuto inficiare sull'efficacia e l'efficienza degli algoritmi, si è scelto di escludere:

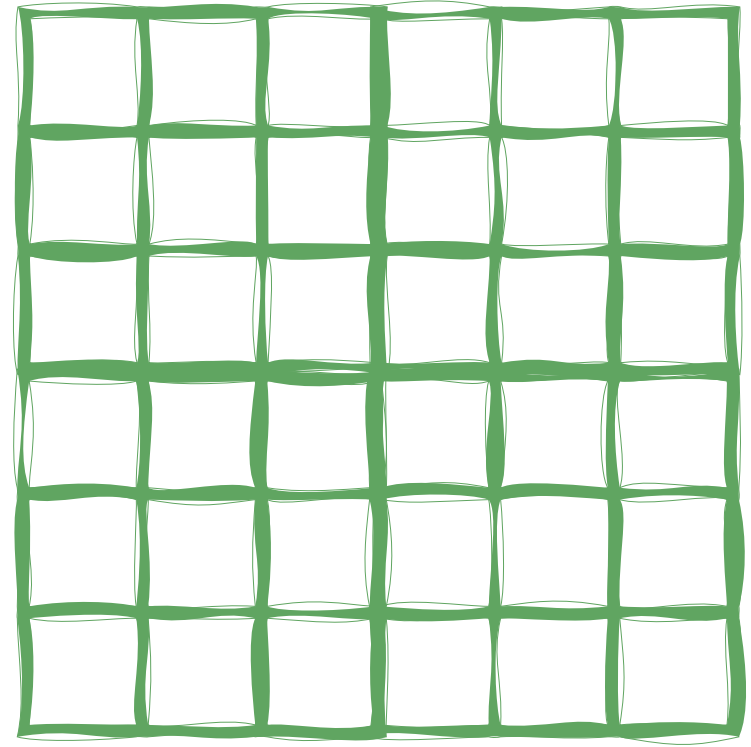
- Termini presenti in **meno di 20 recensioni**
- Termini presenti in **più del 30% delle recensioni**

NMF representation

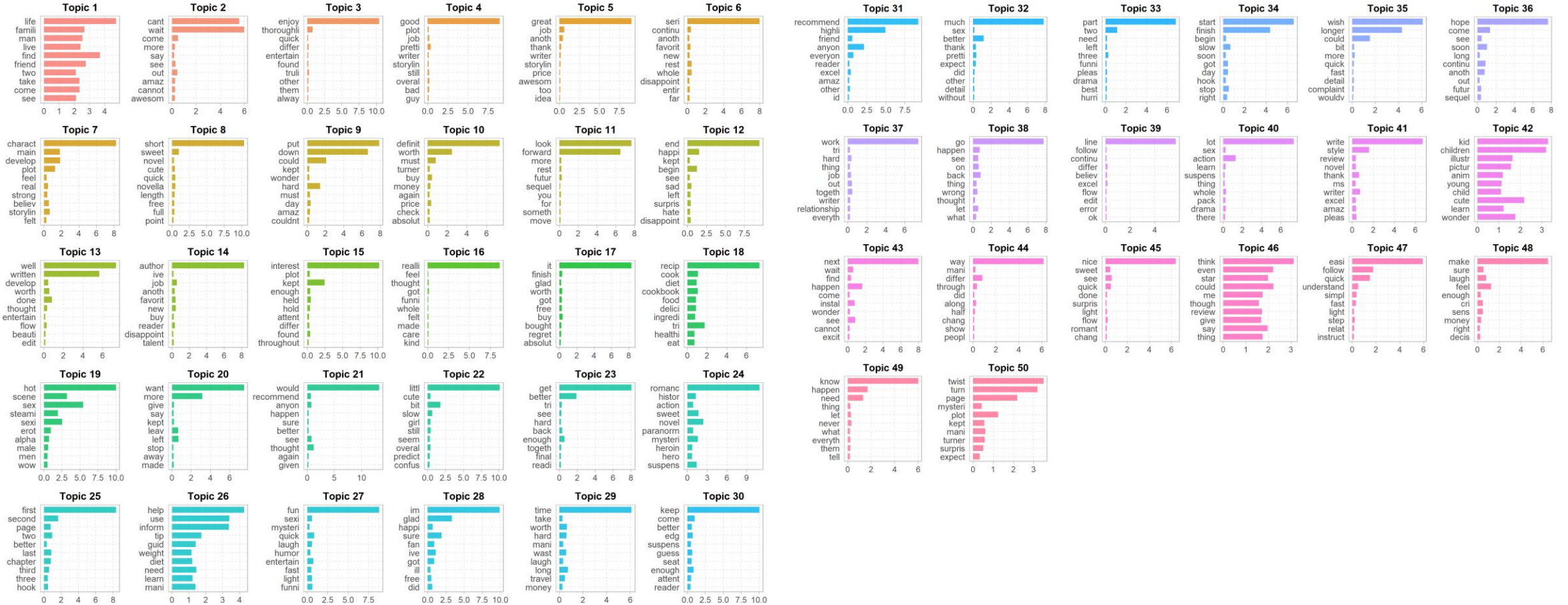
Recensioni Raw



Rappresentazione
TF-IDF con uni-grammi

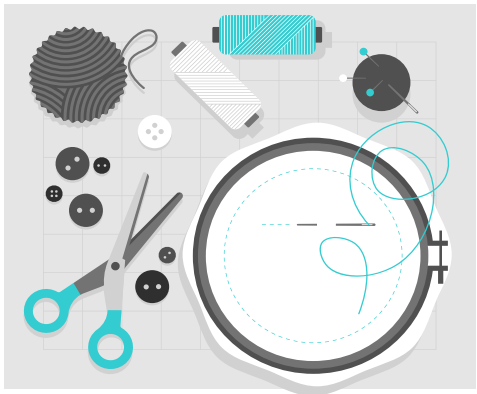


NMF top 50 Topic



Confronto NMF Topic con Book Category

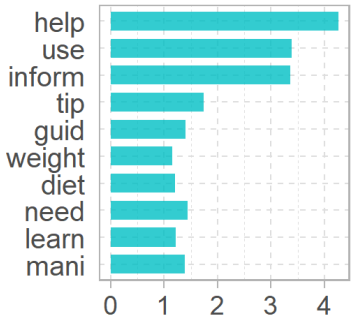
Crafts, Hobbies & Home



Book category



Topic 26

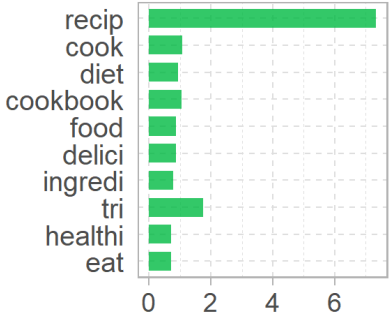


Main Topic

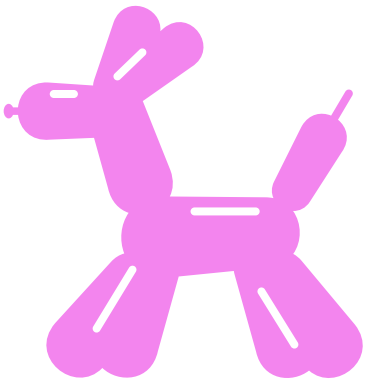
Cookbooks, Food & Wine



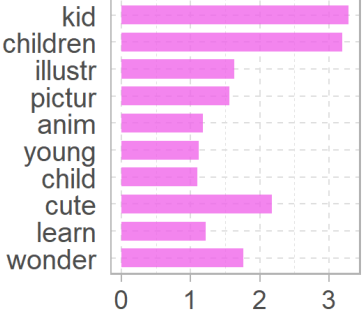
Topic 18



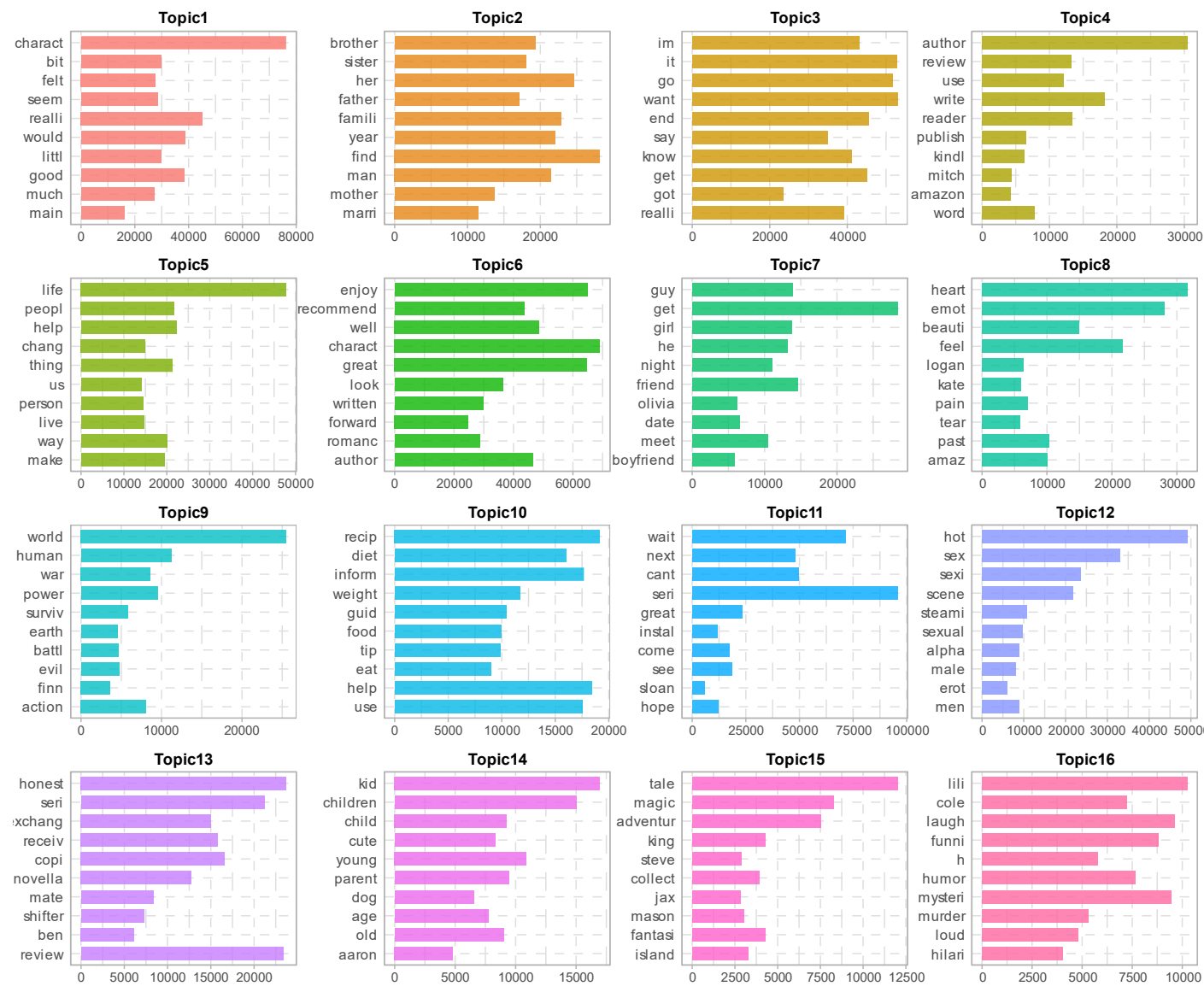
Children's eBooks



Topic 42



LDA Topic



RISULTATI PROGETTO



PERFORMANCE

Il modello classifica in maniera corretta 9 recensioni su 10, indipendentemente dalla loro polarità

SIGNIFICATO

Le parole più discriminanti in un verso e nell'altro emerse dal modello sono logiche e interpretabili

INTERPRETABILITÀ

I topic emersi sia con NMF che con LDA sono interpretabili e alcuni vanno oltre la pura sfera dell'opinione

COERENZA

Il topic NMF sono risultati allineati a quelle che sono le categorie d'appartenenza dei libri

**GRAZIE PER
L'ATTENZIONE**

