



Sentiment Analysis su Amazon Reviews

Scopo del progetto

- Importare le recensioni e modellizzare le entità
- Analizzare e distinguere i topic che caratterizzano un utente
- Distinguere il sentiment dei diversi topic
- Creare una view che permette di visualizzare le analisi fatte

Modellizzazione Dati



Appoggiandosi ad un database **SQLite** e modellizzando i dati con **Django** è stato più semplice gestire le entità rilevate

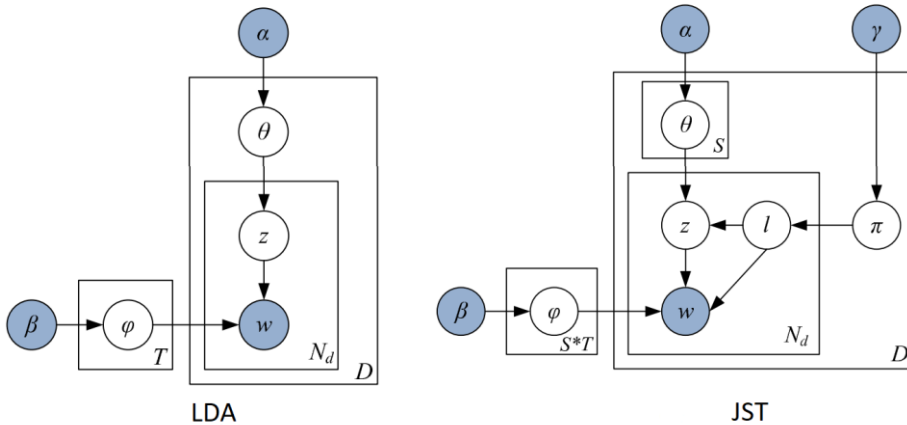
- **User:** (id, num_rating, average_score, variance, experience, pos_words, neg_words)
- **Product:** (id, num_rating, average_score, variance, experience, words)
- **Rating:** (id, productid, userid, score, text, timestamp)



Dataset alternativo

- Formato .tsv
- 85% voti positivi
- Pochissime recensioni per utente
- Assenza timestamp
- Formato .json
- Voti più distribuiti
- Utenti con parecchie recensioni
- Presenza timestamp

Join Sentiment/Topic Model



Modello basato su LDA.

Calcola sia topic che sentiment simultaneamente

- Necessario fare diverse prove per gli iperparametri α , β e γ
- E' importante sapere il numero di topic da trovare
- Selezionare un adeguato numero di parole per topic

Creazione modello

Possibilità di selezionare gli utenti da utilizzare per creare il modello



Amazon reviews analytics

[DATASET](#) [DASHBOARD](#) [USERS](#) [PRODUCTS](#) [REVIEWS](#) [TOPICS](#)

N° words N° topics N° iterations

10

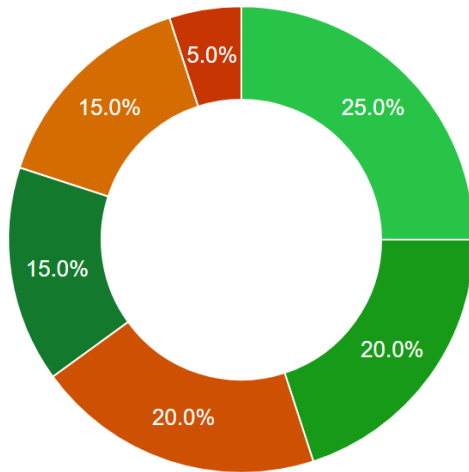
5

20

Make model from selected users

<input type="checkbox"/>	ID ↑ ↓	N° reviews ↑ ↓	Average score ↑ ↓	Variance score ↑ ↓
<input type="checkbox"/>	A00177463W0XWB16A9O0 5	13	4.385	0.544
<input type="checkbox"/>	A022899328A0QROR32DC T	10	3.300	2.610
<input type="checkbox"/>	A04309042SDSL8YX2HRR 7	5	3.800	0.960
<input type="checkbox"/>	A068255029AHTHDXZURN U	9	4.333	0.889
<input type="checkbox"/>	A06944662TFWOKKV4GJK X	9	4.778	0.173
<input type="checkbox"/>	A1004703RC79J9	6	3.667	0.889
<input type="checkbox"/>	A1006HCQDMYC5W	9	4.333	1.111
<input type="checkbox"/>	A1008DPSP6KC9J	7	4.286	0.776
<input type="checkbox"/>	A100DXY4SLAMP	7	5.000	0.000
<input type="checkbox"/>	A100I4UAHGQCF6	5	4.200	2.560
<input type="checkbox"/>	A100L918633LUO	9	4.667	0.889
<input type="checkbox"/>	A100VQNP6I54HS	15	5.000	0.000

Prodotto



■ Topic_1 ■ Topic_3 ■ Topic_5 ■ Topic_0 ■ Topic_2 ■ Topic_4

Distribuzione di topic per prodotto

- Parole che caratterizzano il prodotto calcolate con **LDA**
- Intersezione con i topic per capire la distribuzione

Utente

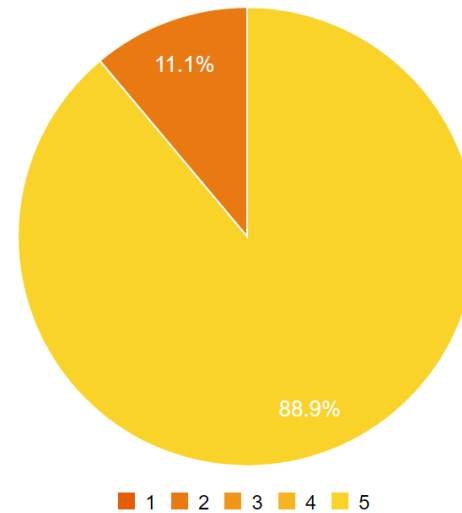
Parole positive (LDA)

flavor wrote differ realli celiac invest similiar textur buyer cream actual cracker
ident everyon belowi better cooki

Parole negative (LDA)

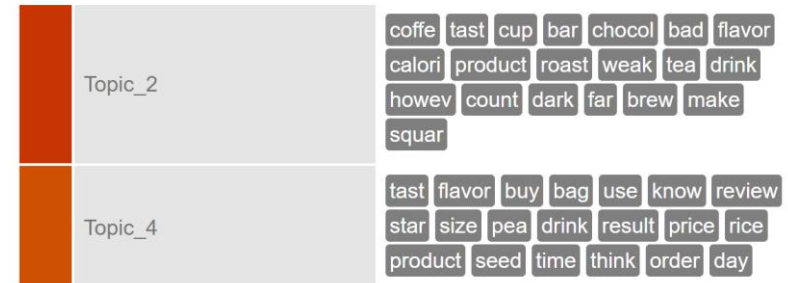
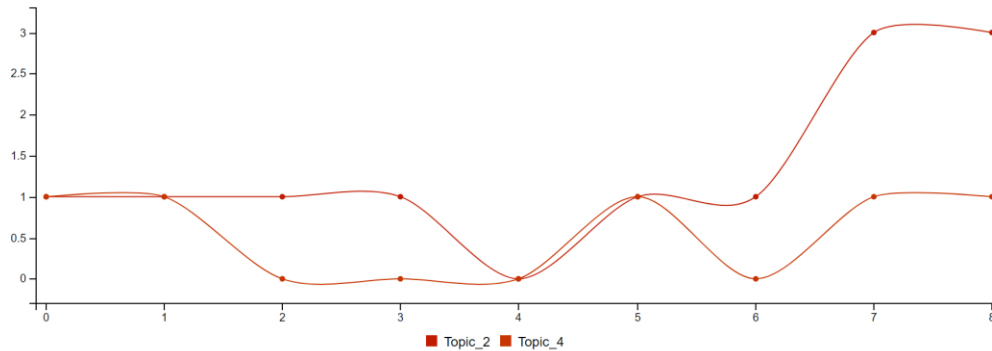
absolut anoth product small textur enfamil refus thing sever disappoint confus truli
review posit

Distribuzione dei commenti per utente:

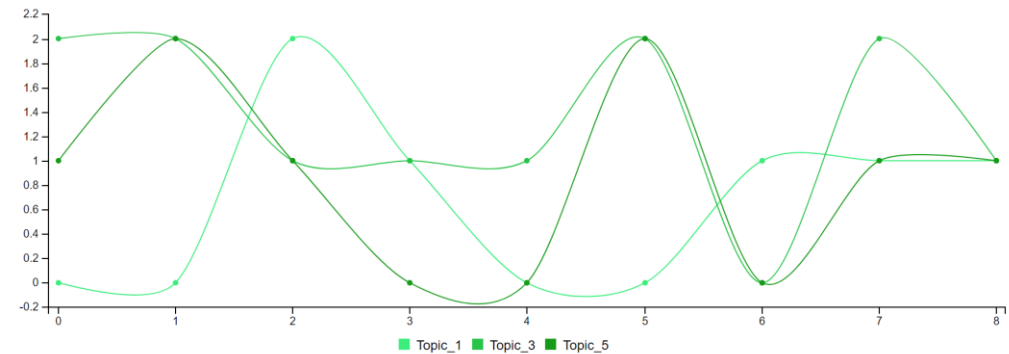


Utente (topic nel tempo)

Negative topics



Positive topics



Conclusione

- Iperparametri fondamentali nella creazione del modello
- Importante selezionare con cura il numero dei topic
- Poche parole (ma anche troppe!) per topic portano il modello a non trovare la giusta distribuzione per utente
- Utilizzare un numero accettabile di iterazioni