

Text Mining and Search

Amazon reviews classification

Obiettivo

Preprocessing

Text representation and Modeling

Conclusioni

Ivan Mera 783086

OBIETTIVO

Generale: Creare una rete neurale per classificare testi

Specifico: Confrontare due tecniche di text representation



Dataset

From



To



Quattro macro categorie:

- Industrial
- Luxury
- Musical Instruments
- Appliances



GloVe

*EDA
&
Feature
Engineering*

*Text
Preprocessing*

Operazioni

- Controllo missing values
- Numero di recensioni per prodotto
- Variabile target



Dataset

From



To



Quattro macro categorie:

- Industrial
- Luxury
- Musical Instruments
- Appliances



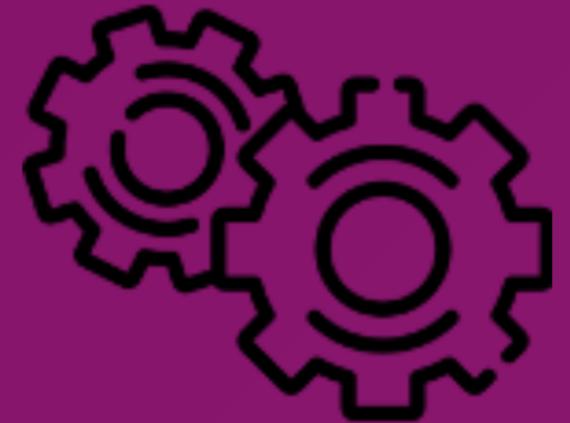
GloVe

*EDA
&
Feature
Engineering*

*Text
Preprocessing*

Operazioni

- Upper case
- Link web
- Contrazioni
- Emoji
- Caratteri speciali
- Spazi bianchi
- Stop words
- Lemmatization



Text representation

Molte possibilità..
Qual è la scelta migliore?



*Text representation
techniques*

Modeling

Results

Approches

- TF-IDF
- WORD EMBEDDINGS



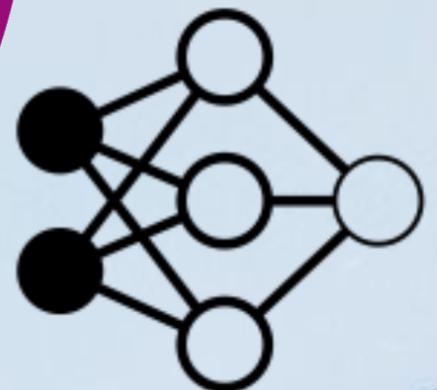
Neural networks architectures

- Dense layers
- CNN + GRU

SUDDIVISIONE DATI:

Train/test: 80% - 20%

Train/Validation: 90% - 10%

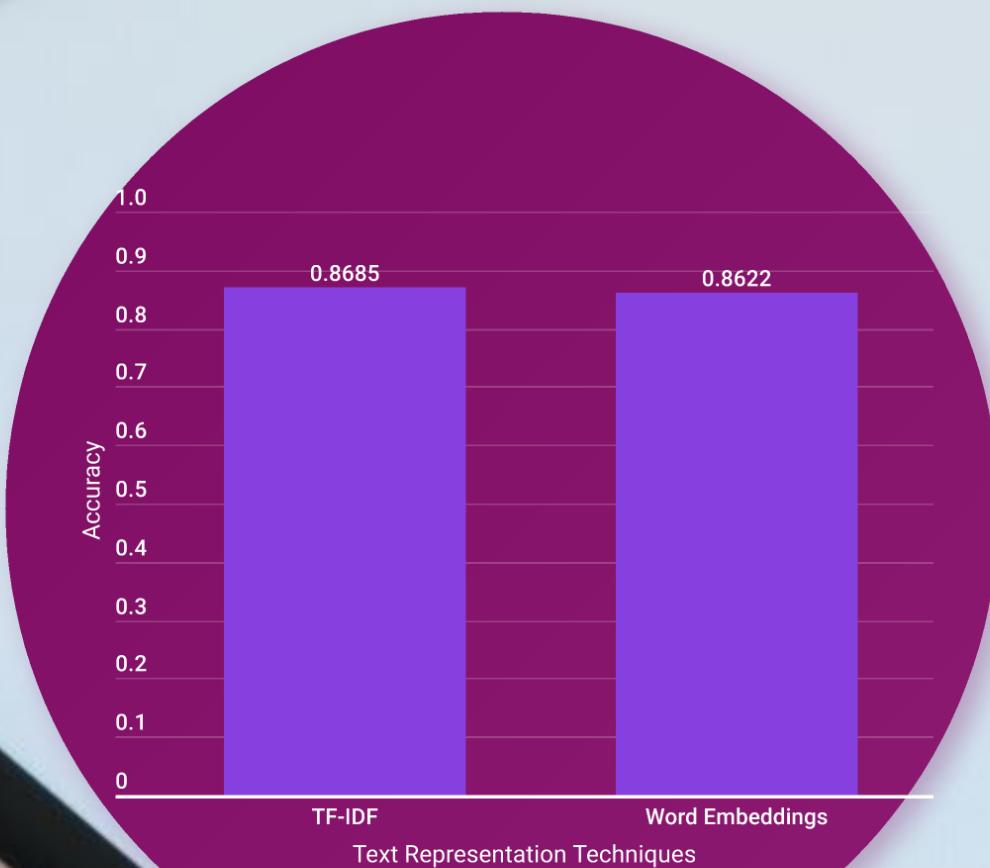




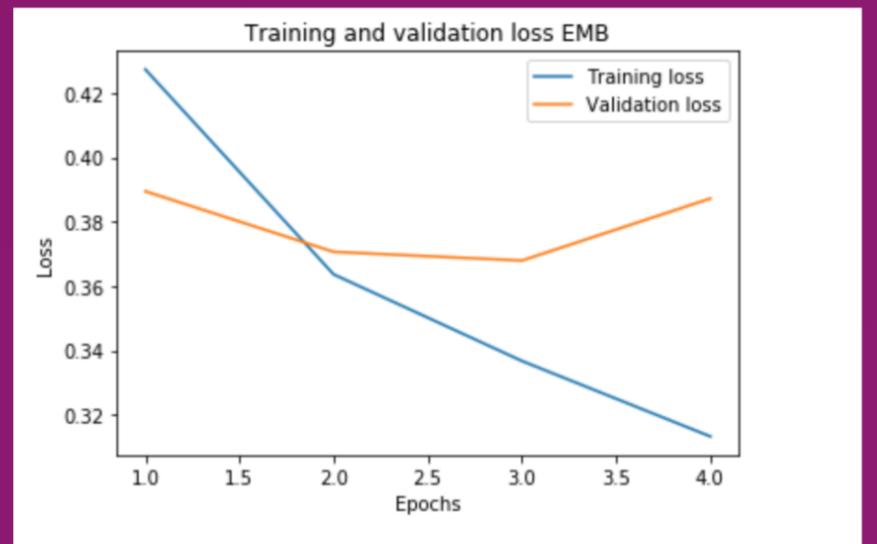
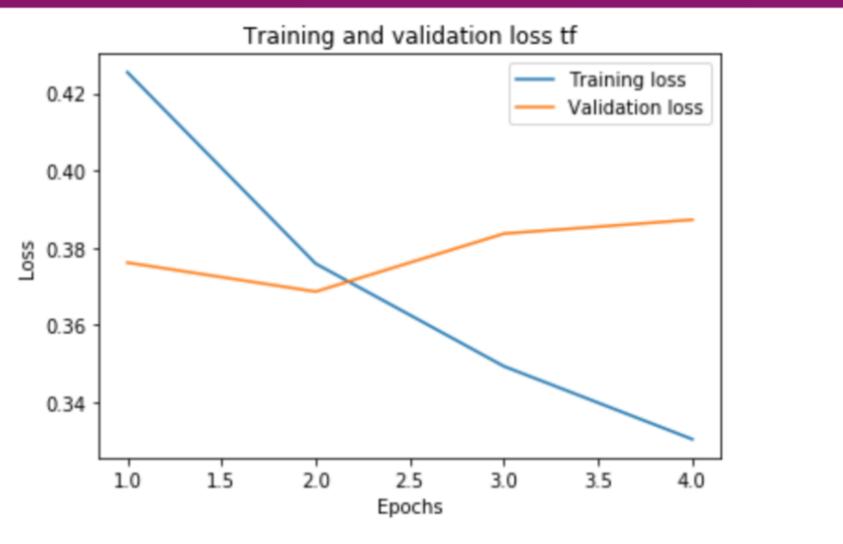
WHO WINS?

Accuracy

Robustness



Overfitting?





IMPROVEMENTS

- Più capacità computazionale
- Più training data
- Ottimizzazione