Università degli Studi di Milano - Bicocca

# The 20 Newsgroups

di Beatrice Fumagalli e Matteo Porcino

Progetto di Text Mining and Search
Appello di gennaio

# Dataset 20—newsgroups

19997 documenti suddivisi in 20 newsgroups:

- alt.atheism
- comp.graphics
- comp.os.ms-windows.misc
- comp.sys.ibm.pc.hardware
- comp.sys.mac.hardware
- comp.windows.x
- misc.forsale
- rec.autos
- rec.motorcycles
- rec.sport.baseball

- rec.sport.hockey
- sci.crypt
- sci.electronics
- sci.med
- sci.space
- soc.religion.christian
- talk.politics.guns
- talk.politics.mideast
- talk.politics.misc
- talk.religion.misc

# Text pre-processing

Dataset suddiviso in *Training set* (70%, 13997 articoli) *e Test set* (30%, 6000 articoli)

- Eliminazione dei metadata
- Eliminazione delle contrazioni
- Rimozione di URLs e indirizzi e-mail
- Tokenization
- Rimozione di righe vuote e tabs
- Rimozione della punteggiatura ed eventuali ulteriori spazi vuoti
- Eliminazione delle stringhe numeriche
- Normalizzazione dei caratteri in minuscolo
- Rimozione delle stop-words
- Eliminazione dei caratteri non unicode
- Eliminazione della punteggiatura
- Stemming e Lemmatization
- Eliminazione dei tokens di lunghezza inferiore a 2

Utilizzo della libreria NLTK

# Text Representation

- **Term Frequencies (TF):** la Term Frequency $tf_{t,d}$ del termine t nel documento d è definito come il numero di volte che t si verifica in d.

- **Term Frequency times Inverse Document Frequency (TF-IDF):** il peso tf-idf di un termine è il prodotto del suo peso tf e del suo peso idf.

Utilizzo della libreria scikit-learn
Ciascuna rappresentazione fornisce una matrice Document-Term di dimensioni 13997x119613

# Text Classification

## Multinomial Naive Bayes

- TF: accuracy = 0.812
- TF-IDF: accuracy = 0.848

## Logistic Regression

- TF: accuracy = 0.809
- TF-IDF: accuracy = 0.823

## Random Forest

- TF: accuracy = 0.619
- TF-IDF: accuracy = 0.626

## Support Vector Machine

- TF: accuracy = 0.754
- TF-IDF: accuracy = 0.840

# Risultati

## Naive Bayes, TF-IDF:

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| alt.atheism | 0.67 | 0.72 | 0.69 | 274 |
| comp.graphics | 0.79 | 0.80 | 0.79 | 301 |
| comp.os.ms-windows.misc | 0.84 | 0.77 | 0.80 | 308 |
| comp.sys.ibm.pc.hardware | 0.72 | 0.85 | 0.78 | 284 |
| comp.sys.mac.hardware | 0.86 | 0.88 | 0.87 | 297 |
| comp.windows.x | 0.88 | 0.89 | 0.89 | 295 |
| misc.forsale | 0.88 | 0.69 | 0.77 | 306 |
| rec.autos | 0.91 | 0.89 | 0.90 | 320 |
| rec.motorcycles | 0.95 | 0.96 | 0.95 | 324 |
| rec.sport.baseball | 0.97 | 0.95 | 0.96 | 309 |
| rec.sport.hockey | 0.95 | 0.98 | 0.97 | 274 |
| sci.crypt | 0.93 | 0.94 | 0.93 | 278 |
| sci.electronics | 0.86 | 0.86 | 0.86 | 307 |
| sci.med | 0.93 | 0.92 | 0.93 | 297 |
| sci.space | 0.90 | 0.92 | 0.91 | 287 |
| soc.religion.christian | 0.89 | 0.92 | 0.90 | 308 |
| talk.politics.guns | 0.82 | 0.88 | 0.85 | 296 |
| talk.politics.mideast | 0.89 | 0.94 | 0.91 | 329 |
| talk.politics.misc | 0.74 | 0.73 | 0.74 | 319 |
| talk.religion.misc | 0.55 | 0.46 | 0.50 | 287 |
| | | | | |
| micro avg | 0.85 | 0.85 | 0.85 | 6000 |
| macro avg | 0.85 | 0.85 | 0.85 | 6000 |
| weighted avg | 0.85 | 0.85 | 0.85 | 6000 |

# Risultati

Support Vector Machine, TF-IDF:



| | precision | recall | f1-score | support |
|---|---|---|---|---|
| alt.atheism | 0.70 | 0.71 | 0.70 | 274 |
| comp.graphics | 0.82 | 0.80 | 0.81 | 301 |
| comp.os.ms-windows.misc | 0.80 | 0.81 | 0.80 | 308 |
| comp.sys.ibm.pc.hardware | 0.74 | 0.79 | 0.76 | 284 |
| comp.sys.mac.hardware | 0.86 | 0.87 | 0.87 | 297 |
| comp.windows.x | 0.86 | 0.87 | 0.86 | 295 |
| misc.forsale | 0.77 | 0.81 | 0.79 | 306 |
| rec.autos | 0.91 | 0.86 | 0.88 | 320 |
| rec.motorcycles | 0.95 | 0.95 | 0.95 | 324 |
| rec.sport.baseball | 0.97 | 0.94 | 0.96 | 309 |
| rec.sport.hockey | 0.96 | 0.97 | 0.96 | 274 |
| sci.crypt | 0.94 | 0.94 | 0.94 | 278 |
| sci.electronics | 0.82 | 0.83 | 0.82 | 307 |
| sci.med | 0.92 | 0.92 | 0.92 | 297 |
| sci.space | 0.95 | 0.92 | 0.93 | 287 |
| soc.religion.christian | 0.87 | 0.89 | 0.88 | 308 |
| talk.politics.guns | 0.83 | 0.84 | 0.83 | 296 |
| talk.politics.mideast | 0.91 | 0.91 | 0.91 | 329 |
| talk.politics.misc | 0.73 | 0.69 | 0.71 | 319 |
| talk.religion.misc | 0.49 | 0.48 | 0.49 | 287 |
| | | | | |
| micro avg | 0.84 | 0.84 | 0.84 | 6000 |
| macro avg | 0.84 | 0.84 | 0.84 | 6000 |
| weighted avg | 0.84 | 0.84 | 0.84 | 6000 |