

Introduzione alla Sentiment Analysis - II

Rachele Sprugnoli – rachele.sprugnoli@unicatt.it

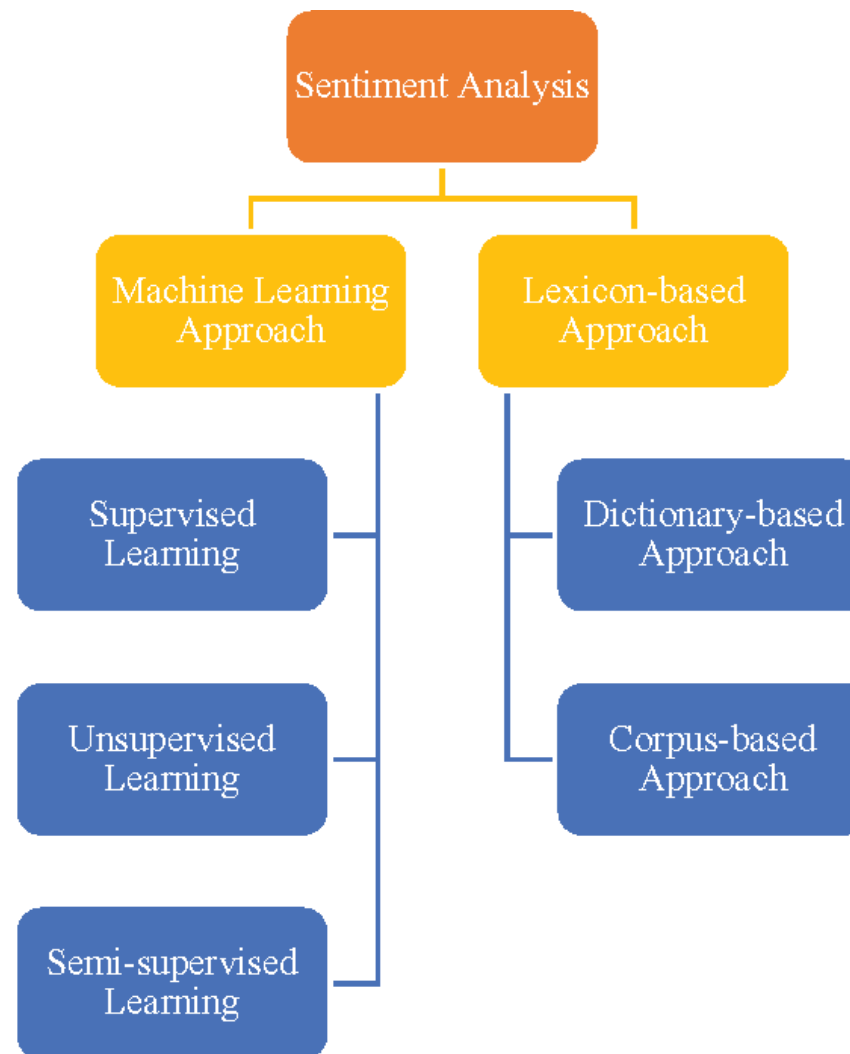
Centro Interdisciplinare di Ricerche per la Computerizzazione
dei Segni dell'Espressione (CIRCSE)



UNIVERSITÀ
CATTOLICA
del Sacro Cuore



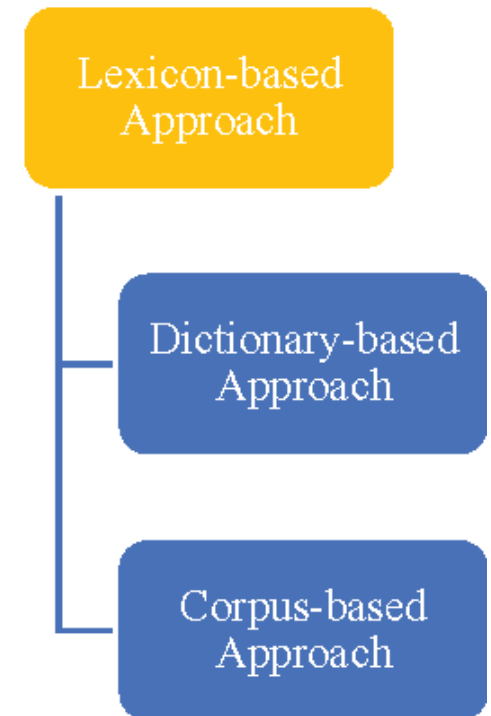
COME SI RISOLVE UN TASK DI SENTIMENT ANALYSIS



APPROCCI BASATI SUL LESSICO

- Basati sull'intuizione che la polarità di testo può essere ottenuta sulla base della polarità delle parole che lo compongono
- Si basano su liste di token, espressioni o lemmi («*sentiment/polarity lexicon*»): ad ogni token o lemma sono associate delle categorie o dei valori numerici che ne quantificano la polarità
Esempi:

- *negative, neutral, positive*
- *very negative, negative, neutral, positive, very positive*
- *-1, 0, +1*
- *-2, -1, 0, +1, +2*
- *valori decimali tra 0 e 1*



COME SI CREANO LESSICI DI SENTIMENT

- **DICTIONARY-BASED:** si usano una lista di termini con associato il sentiment e dei dizionari in formato digitale in cui si cercano sinonimi e contrari di quei termini

- *paura* = **NEG**

SINONIMI = *spavento, terrore* = **NEG**

CONTRARI = *coraggio, audacia* = **POS**

- **CORPUS-BASED:** si usano una lista di termini con associato il sentiment e delle collezioni di testi. Altri termini vengono aggiunti alla lista cercando contesti comuni o strutture sintattiche particolari

La macchina è bella E spaziosa

La macchina è bella MA difficile da guidare

Lexicon-based
Approach

Dictionary-based
Approach

Corpus-based
Approach



LESSICI DI SENTIMENT PER L'ITALIANO - 1

- **OPENER:** più di 25.000 lemmi ed espressioni con PoS, 877 con sentiment assegnato manualmente, gli altri in maniera automatica

```
<LexicalEntry id="id_24468" partOfSpeech="noun">  
  <Lemma writtenForm="amico"/>  
  <Sense>  
    <Confidence score="0.5" method="automatic"/>  
    <Sentiment polarity="positive"/>  
    <Domain/>  
  </Sense>  
</LexicalEntry>
```

<https://github.com/opener-project/VU-sentiment-lexicon/tree/master/VUSentimentLexicon/IT-lexicon>



LESSICI DI SENTIMENT PER L'ITALIANO - 2

- **Sentix (Sentiment Italian Lexicon)**: più di 74.000 token ed espressioni con PoS, sentiment assegnato tramite allineamento di risorse pre-esistenti
 - sentiment diverso in base al significato

PAROLA	PoS	senso	POS_score	NEG_score	polarity	intensity
amico	n	09785042	0.25	0	1.0	0.25
amico	n	09877951	0.125	0	1.0	0.125
amico	n	10112591	0.125	0	1.0	0.125
amico	n	10686073	0.375	0	1.0	0.375

“someone who shares your feelings or opinions and hopes that you will be successful”



LESSICI DI SENTIMENT PER L'ITALIANO - 3

- **Distributional Polarity Lexicon:** > 75.000 lemmi con PoS, sentiment associato automaticamente a partire da una collezione di tweet

parola::PoS positività, negatività, neutralità

amico::s 0.44531634,0.24274117,0.31194243

amico::a 0.7146159,0.15039273,0.1349914

#sanremo2012::h 0.26471627,0.28195098,0.45333272

#buonavita::h 0.71994877,0.052290898,0.22776033



LESSICI DI SENTIMENT PER L'ITALIANO - 4

- **SenticNet**: più di 23.000 lemmi ed espressioni rappresentanti concetti di senso comune con associati vari valori di polarità ed emozione

```
▼<rdf:RDF xmlns:rdf="http://w3.org/1999/02/22-rdf-syntax-ns#">
  ▼<rdf:Description rdf:about="http://sentic.net/api/it/concept/amico">
    <rdf:type rdf:resource="http://sentic.net/api/concept"/>
    <text xmlns="http://sentic.net">amico</text>
    ▼<semantics xmlns="http://sentic.net">
      <concept xmlns="http://sentic.net" rdf:resource="http://sentic.net/api/it/concept/amicizia"/>
      <concept xmlns="http://sentic.net" rdf:resource="http://sentic.net/api/it/concept/grande"/>
      <concept xmlns="http://sentic.net" rdf:resource="http://sentic.net/api/it/concept/convivente"/>
      <concept xmlns="http://sentic.net" rdf:resource="http://sentic.net/api/it/concept/si_basano_su_altra_persona"/>
      <concept xmlns="http://sentic.net" rdf:resource="http://sentic.net/api/it/concept/divertirsi_insieme"/>
    </semantics>
    ▼<sentic xmlns="http://sentic.net">
      <pleasantness xmlns="http://sentic.net" rdf:datatype="http://w3.org/2001/XMLSchema#float">0.72</pleasantness>
      <attention xmlns="http://sentic.net" rdf:datatype="http://w3.org/2001/XMLSchema#float">0</attention>
      <sensitivity xmlns="http://sentic.net" rdf:datatype="http://w3.org/2001/XMLSchema#float">0</sensitivity>
      <aptitude xmlns="http://sentic.net" rdf:datatype="http://w3.org/2001/XMLSchema#float">0.879</aptitude>
    </sentic>
    ▼<moodtags xmlns="http://sentic.net">
      <concept xmlns="http://sentic.net" rdf:resource="http://sentic.net/api/it/concept/felicità"/>
      <concept xmlns="http://sentic.net" rdf:resource="http://sentic.net/api/it/concept/ammirazione"/>
    </moodtags>
    ▼<polarity xmlns="http://sentic.net">
      <value xmlns="http://sentic.net">positive</value>
      <intensity xmlns="http://sentic.net" rdf:datatype="http://w3.org/2001/XMLSchema#float">0.533</intensity>
    </polarity>
  </rdf:Description>
</rdf:RDF>
```



LESSICI DI SENTIMENT PER L'ITALIANO - 5

- **NRC-Emotion-Lexicon:** quasi 14.000 lemmi ed espressioni tradotti automaticamente dall'inglese, assegnazione di positività o negatività + assegnazione delle emozioni di Plutchik basata sull'annotazione manuale di vari annotatori non esperti

English	Italian	Pos	Neg	Anger	Anticipation	Disgust	Fear	Joy	Sadness	Surprise	Trust
friend	amico	1	0	0	0	0	0	1	0	0	1
nausea	nausea	0	1	0	0	1	0	0	0	0	0
unnatural	innaturale	0	1	0	0	1	1	0	0	0	0

<https://saifmohammad.com/WebPages/NRC-Emotion-Lexicon.htm>



LESSICI DI SENTIMENT PER L'ITALIANO - 6

- **NRC-VAD-Lexicon**: più di 19.000 token ed espressioni tradotti automaticamente dall'inglese, assegnazione dei valori di Valence, Arousal e Dominance basata sull'annotazione manuale di vari annotatori non esperti

Word	Italian-it	Valence	Arousal	Dominance
friend	amico	0.906	0.413	0.573
pal	amico	0.812	0.426	0.596
nausea	nausea	0.104	0.781	0.273
natural	naturale	0.854	0.118	0.481
unnatural	innaturale	0.153	0.587	0.434

<https://saifmohammad.com/WebPages/nrc-vad.html>



LESSICI DI SENTIMENT PER L'ITALIANO - 7

- **Depeche Mood**: più di 82.000 parole ed espressioni con o senza PoS, assegnazione automatica di 5 emozioni derivata dai commenti online agli articoli del Corriere della Sera

TOKEN#PoS	INDIGNATO	PREOCCUPATO	TRISTE	DIVERTITO	SODDISFATTO	freq
raccoglimento#n	0.078	0.005	0.829	0.049	0.036	4
sims#n	0.0	0.0	0.0	0.687	0.312	15
fiabesco#a	0.012	0.0	0.028	0.0	0.959	3
squilibri#n	0.105	0.894	0.0	0.0	0.0	1
euroburocrati#n	0.972	0.0	0.0	0.0	0.027	3

<https://github.com/marcoguerini/DepecheMood/releases/tag/v2.0>



DATI ANNOTATI PER L'ITALIANO - 1

- **Sentipolc:** Sentiment Polarity Classification
 - Corpus di tweet annotati su quattro dimensioni:
 1. Soggettività: soggettivo versus oggettivo
 2. Polarità: negativo, positivo, mixed
 3. Ironia: ironico, non ironico
 4. Presenza di linguaggio figurato: polarità letterale, polarità non letterale

subj	opos	oneg	iro	lpos	lneg	description and explanatory tweet in Italian
0	0	0	0	0	0	objective <i>l'articolo di Roberto Ciccarelli dal manifesto di oggi</i> http://fb.me/1BQVy5Wak
1	0	0	0	0	0	subjective with neutral polarity and no irony <i>Primo passaggio alla #strabrollo ma secondo me non era un iscritto</i>
1	1	0	0	1	0	subjective with positive polarity and no irony <i>splendida foto di Fabrizio, pluri cliccata nei siti internazionali di Photo Natura</i> http://t.co/GWoZqbxAuS
1	0	1	0	0	1	subjective with negative polarity and no irony <i>Monti, ripensaci: l'inutile Torino-Lione inguaia l'Italia: Tav, appello a Mario Monti da Mercalli, Cicconi, Pont...</i> http://t.co/3CazKS7Y
1	1	1	0	1	1	subjective with both positive and negative polarity (mixed polarity) and no irony <i>Dati negativi da Confindustria che spera nel nuovo governo Monti. Castiglione: "Avanti con le riforme"</i> http://t.co/kIKnbFY7

<http://www.di.unito.it/~tutreeb/sentipolc-evalita16/data.html>



DATI ANNOTATI PER L'ITALIANO - 2

- **ABSITA:** Aspect-based Sentiment Analysis
 - Corpus di frasi tratte da recensioni di hotel di Booking annotati considerando 8 caratteristiche (e.g. pulizia, location) su due dimensioni:
 1. Presenza o meno della caratteristica
 2. Sentiment positivo o negativo per ciascuna caratteristica

cleanliness_ presence	cleanliness_ positive	cleanliness_ negative	location_ presence	location_ positive	location_ negative	sentence
0	0	0	1	1	0	La posizione, davvero invidiabile
0	0	0	1	0	1	Un po' fuori mano
1	1	0	1	1	0	Albergo pulito, in un'ottima posizione
1	0	1	0	0	0	Sporco a livelli colossali



DATI ANNOTATI PER L'ITALIANO - 3

- **italian-sentiment-analysis**: Aspect-based Sentiment Analysis
 - Corpus di recensioni di smartphone

schermo[+],touch[+],volume[+],prezzo[+],ram[-]##Come si fa ad creare un cellulare simile : schermo bellissimo , touch bellissimo , volume bellissimo , prezzo incredibile , 2 alloggi per SIM Ma con così poca RAM Solo 512Mega

batteria[+],design[+],schermo[+]##Batteria impressionante, bel design e schermo piacevole.



DATI ANNOTATI PER L'ITALIANO - 4

- **HaSpeeDe:** Hate Speech Detection
 - Corpus di tweet e post di Facebook annotati per la presenza di discorsi d'odio

Fonte	Tweet/Post	HateSpeech
Facebook	ITALIANI. BRUCIAMO LA STREGA	1
Facebook	Io mi auguro che gli italiani aprano finalmente gli occhi. @magdicristiano @angelo_ra_ Ci vuole la guerra per salvare	0
Twitter	l'Italia dai criminali filo islamici. https://t.co/TqHWKr74jB I traffici di organi seguono le rotte dei migranti. Summit @CasinaPioIV in Vaticano con medici giuristi politici da...	1
Twitter	https://t.co/kBGuhoHtpL	0



DATI ANNOTATI PER L'ITALIANO - 5

- **Italian Twitter Corpus of Hate Speech:** Hate Speech Detection

- Corpus di tweet annotati su 7 dimensioni:

1. Presenza di hate speech = yes
2. Target del tweet = religion
3. Grado di aggressività = strong
4. Grado di offensività = strong
5. Presenza di ironia = no
6. Presenza di stereotipi = yes
7. Grado di intensità = 3



Replying to @1282Ugo and @sputnik_italia

Bene,ricomincia la distruzione dei Porci terroristi

[Translate Tweet](#)

2:37 PM · Apr 12, 2017 · [Twitter for Android](#)

<https://github.com/msang/hate-speech-corpus>



DATI ANNOTATI PER L'ITALIANO - 6

- **Happy Parents:** Sentiment Polarity Classification
 - Corpus di tweet sul tema della fecondità e genitorialità annotati su 3 dimensioni:
 1. Soggettività
 2. Polarità con o senza ironia
 3. Argomento

Tweet	Annotazione	Argomento
@criharry bimba mia ti voglio tanto bene	POS	Being parent
Nella mia famiglia non parliamo tanto, però alla fine ci vogliamo bene!	MIXED	Daily live
Il trionfo di Renzi indebolisce vita e famiglia http://t.co/yAd1l10owN	NEG	Fertility & politics
Non ero preparato a fare il figlio, figuriamoci il genitore.	NEG	Becoming parents
Figliolo, lo vedi tutto questo disagio? Un giorno sarà tuo.	HUMNEG	Being parent



DATI ANNOTATI PER L'ITALIANO - 7

- **WhatsApp Dataset:** Cyberbullying Detection
 - Corpus di chat di WhatsApp raccolte tramite simulazione in classi seconde di scuole medie annotate in base a due dimensioni:
 1. Ruolo di chi scrive: e.g. vittima, bullo...
 2. Tipo di cyberbullismo: e.g. minacce, body shame, sessismo

Che fallito, vatti a chiudere nella fogna, finocchio!

[che fallito,]_{Insult: General Insult} [vatti a chiudere nella fogna,]_{Curse or Exclusion}

[finocchio!]_{Insult-Discrimination: Sexism}



QUANTO BENE FUNZIONANO I SISTEMI?

SISTEMI MACHINE LEARNING SULL'ITALIANO

- **Subjectivity Detection:** 0,74 F-score (2016)
- **Sentiment Classification** (Twitter): 0,66 F-score (2016)
- **Aspect-based Sentiment Analysis:** 0,81 F-score (2018)
- **Irony Detection:** 0,73 F-score (2018)
- **Irony Classification:** 0,52 F-score (2018)
- **Misogeny Detection:** 0,84 Accuracy (2018)
- **Misogeny Classification:** 0,41 Accuracy (2018)
- **Hate Speech Detection** (Twitter): 0,80 F-score (2018)
- **Hate Speech Detection** (FaceBook): 0,83 F-score (2018)

Performance dei migliori sistemi di EVALITA 2016 e 2018

