

Andrea Ceron
Luigi Curini
Stefano M. Iacus

Sei

Social Media e Sentiment Analysis

L'evoluzione dei fenomeni sociali
attraverso la Rete



Springer

Sxi – Springer per l’Innovazione

Sxi – Springer for Innovation

Volume 9

For further volumes:
<http://www.springer.com/series/10062>

Andrea Ceron • Luigi Curini • Stefano M. Iacus

Social Media e Sentiment Analysis

L'evoluzione dei fenomeni sociali
attraverso la Rete

 Springer

Andrea Ceron
Dipartimento di Scienze Sociali
e Politiche
Università degli Studi di Milano
Milano, Italia

Luigi Curini
Dipartimento di Scienze Sociali
e Politiche
Università degli Studi di Milano
Milano, Italia

Stefano M. Iacus
Dipartimento di Economia,
Management e Metodi Quantitativi
Università degli Studi di Milano
Milano, Italia

Sxi – Springer per l’Innovazione / Sxi – Springer for Innovation

ISSN: 2239-2688

ISSN: 2239-2696 (electronic)

ISBN 978-88-470-5531-5

ISBN 978-88-470-5532-2 (eBook)

DOI 10.1007/978-88-470-5532-2

Springer Milan Heidelberg New York Dordrecht London

© Springer-Verlag Italia 2014

Quest’opera è protetta dalla legge sul diritto d’autore e la sua riproduzione anche parziale è ammessa esclusivamente nei limiti della stessa. Tutti i diritti, in particolare i diritti di traduzione, ristampa, riutilizzo di illustrazioni, recitazione, trasmissione radiotelevisiva, riproduzione su microfilm o altri supporti, inclusione in database o software, adattamento elettronico, o con altri mezzi oggi conosciuti o sviluppati in futuro, rimangono riservati. Sono esclusi brevi stralci utilizzati a fini didattici e materiale fornito ad uso esclusivo dell’acquirente dell’opera per utilizzazione su computer. I permessi di riproduzione devono essere autorizzati da Springer e possono essere richiesti attraverso RightsLink (Copyright Clearance Center). La violazione delle norme comporta le sanzioni previste dalla legge.

Le fotocopie per uso personale possono essere effettuate nei limiti del 15% di ciascun volume dietro pagamento alla SIAE del compenso previsto dalla legge, mentre quelle per finalità di carattere professionale, economico o commerciale possono essere effettuate a seguito di specifica autorizzazione rilasciata da CLEARedi, Centro Licenze e Autorizzazioni per le Riproduzioni Editoriali, e-mail autorizzazioni@clearedi.org e sito web www.clearedi.org.

L’utilizzo in questa pubblicazione di denominazioni generiche, nomi commerciali, marchi registrati, ecc., anche se non specificatamente identificati, non implica che tali denominazioni o marchi non siano protetti dalle relative leggi e regolamenti.

Le informazioni contenute nel libro sono da ritenersi veritiere ed esatte al momento della pubblicazione; tuttavia, gli autori, i curatori e l’editore declinano ogni responsabilità legale per qualsiasi involontario errore od omissione. L’editore non può quindi fornire alcuna garanzia circa i contenuti dell’opera.

9 8 7 6 5 4 3 2 1

Layout copertina: Beatrice E, Milano

Impaginazione: PTP-Berlin, Protago TEX-Production GmbH, Germany (www.ptp-berlin.eu)

Stampa: Grafiche Porpora, Segrate (MI)

Springer-Verlag Italia S.r.l., Via Decembrio 28, I-20137 Milano

Springer-Verlag fa parte di Springer Science+Business Media (www.springer.com)

Presentazioni

Se proprio non ci fossero altri buoni motivi per leggere questo volume, mi basterebbe il fatto che vi sono citati due tra i romanzi con cui sono cresciuto: la *Trilogia della Fondazione*¹ di Isaac Asimov (che aveva immaginato una scienza, la psicostoria, in grado di prevedere il futuro a patto di poter esaminare il comportamento di un numero sufficientemente grande di persone), e *Guida galattica per gli autostoppisti* di Douglas Adams. Ci fosse stato anche il *Doctor Who*, ne avrei proposto l'adozione nella scuola dell'obbligo.

In realtà (e per fortuna per voi che lo avete in mano) di ragioni di lettura ve ne sono molte. Anche se la politica e l'informazione italiana per la maggior parte non se ne sono accorte, viviamo in un'epoca in cui la distinzione tra reale e virtuale non esiste più. L'utilizzo della Rete è abitudine per metà del paese, mentre sono circa 17 milioni i connazionali che ogni giorno lavorano, si innamorano, si relazionano, concludono transazioni commerciali attraverso i social network, primo tra tutti Facebook. Siamo circondati da dispositivi connessi e geolocalizzati: gli smartphone, i portatili, i tablet, le automobili, i bancomat, gli autobus, le videocamere di sorveglianza, i caselli autostradali, gli oggetti di vita quotidiana. Una tendenza destinata crescere in modo esponenziale: tra sette anni vi saranno dieci dispositivi collegati a internet per ogni individuo del pianeta². Il digitale dunque non è *second life*, così si chiamava un ambiente di realtà virtuale che ebbe qualche fortuna anni fa, ma vita vera. È ormai intrecciato in modo indissolubile con le nostre abitudini, alle quali aggiunge nuove dimensioni di profondità.

¹ In realtà Asimov non solo si appassionò alla previsione della storia, ma si applicò anche alla ricerca sociale in politica. In un racconto del 1955, "Franchise" (<http://www.scribd.com/doc/23542910/Asimov-Isaac-The-Complete-Stories-Volume-1>), in italiano "Diritto di voto", arriva a immaginare un sistema di previsione delle intenzioni elettorali talmente evoluto da essere in grado di selezionare un solo individuo, tra tutti gli americani, che riassume in sé caratteri, personalità, volontà dell'intero elettorato, e dunque alle presidenziali vota per tutti. Curiosamente, lo scrittore di fantascienza situa la storia nel novembre del 2008, ovvero in corrispondenza con le elezioni presidenziali che hanno consacrato Barack Obama, al termine della prima campagna svoltasi in modo massiccio anche attraverso internet.

² Morgan Stanley prevede che nel 2020 ci saranno oltre 75 miliardi di dispositivi connessi a internet, mentre la popolazione stimata sarà di 8 miliardi di persone. <http://www.businessinsider.com/75-billion-devices-will-be-connected-to-the-internet-by-2020-2013-10>.

Come mi raccontava Nathan Jurgenson³ di recente, “il web, il digitale, gli smart-phone, Facebook e tutto il resto sono fortemente materiali e fisici [...]. La finzione dell’online come virtuale serve solo a contrapporgli un presunto reale “naturale”, come ideale di vita vera, disconnessa. Un’ideale irraggiungibile, conservatore e che in ultima analisi serve solo a disumanizzare quanti sono immersi nella dimensione aumentata della realtà contemporanea”.

I dati sono il petrolio del XXI secolo. Ciò porta con sé nuovi problemi, come nota Evgeny Morozov⁴ riferendosi alla presunta gratuità di molti servizi del social web, dallo spazio disco alla posta elettronica, dalla videocomunicazione ai servizi di condivisione di video e foto: “Se una prestazione che ha un costo ti viene fornita gratis, allora è meglio che ti preoccupi: significa che la merce sei tu”. Temi che nei prossimi anni porteranno alla necessità di ridefinire criteri di consumo consapevole per il digitale, come accadde nel XX secolo quando, per far fronte all’inquinamento, furono introdotti comportamenti di sostenibilità ambientale oggi scontati. A ciò si aggiungono allarmi per la privacy, sia per la sindrome del Grande Fratello cui lo scandalo Nsa-Prism ha conferito nuova concretezza, sia per i fratellini che ogni giorno utilizzano i nostri dati a nostra insaputa.

Tuttavia, il monitoraggio della Rete in modo aggregato attraverso i big data forniti dai social network fornisce un’opportunità straordinaria alla ricerca sociale: il *nowcasting*,⁵ in altre parole la rilevazione del presente, e la previsione del futuro. Quale sia il valore dei dati raccolti lo comprendono a pieno gli stessi social network: Twitter, ad esempio, consente pieno accesso ai propri dati solo a un club esclusivo di quattro aziende, che rivendono informazioni e servizi a caro prezzo. Ha osservato Carola Frediani: “Chi può salvare i contenuti dei tweet non ha in mano solo un modello di business, ma anche una nuova forma di potere”.⁶ Come leggerete nei capitoli che seguono, attraverso l’analisi dei tweet si cercano oggi di prevedere i fenomeni più diversi: l’atteggiamento dei consumatori verso le aziende, il diffondersi delle malattie, la vittoria a Sanremo e nei talent televisivi, l’andamento dei mercati finanziari, il risultato delle elezioni, la felicità degli italiani.

Va detto che si tratta di scienza nella primissima infanzia e che molti studiosi ritengono non vi siano ancora certezze sufficienti. A questo proposito valga citare tra tutti il paper accademico di Daniel Gayo-Avello, dell’Università di Oviedo, dall’eloquente titolo “I Wanted to Predict Elections with Twitter and all I got was this Lousy Paper”.⁷ Richard Rogers,⁸ al quale poche settimane prima dell’uscita di questo volume avevo chiesto un parere in proposito, mi ha risposto di ritenere che

³ Sociologo, teorico dei social media, contributing editor di “New Inquiry” (<http://nathanjurgenson.tumblr.com/>).

⁴ Per una critica ai social network e al loro impianto economico vedi anche Andrew Keen (<http://ajkeen.com/>) e Michel Bauwens (http://p2pfoundation.net/MichelBauwens/Full_Bio).

⁵ A questo proposito vedi Choi, Varian (<http://www.frbsf.org/economic-research/files/Varian-part.1.pdf>).

⁶ In “Twitter, rivendite autorizzate”, Wired Italia, settembre 2013, pp. 100 e ss.

⁷ Il titolo fa il verso alle magliette souvenir con la scritta “I miei genitori sono stati in vacanza a ... e tutto quel che mi hanno portato è questa maglietta schifosa”. <http://arxiv.org/pdf/1204.6441.pdf>.

⁸ Cattedratico di Nuovi Media e Cultura Digitale all’Università di Amsterdam <http://www.uva.nl/overde-uva/organisatie/medewerkers/content/r/o/r.a.rogers/r.a.rogers.html>.

esistano ancora due problemi significativi nell'uso di Twitter per i sondaggi politici: "Che cosa si possa veramente intendere su Twitter per voto, e quali possano essere i termini reali di paragone delle nostre osservazioni".

In questo libro il merito di Ceron, Curini e Iacus è di non nascondere le difficoltà di interpretare i sentiment attraverso i tweet, ma di cercare di spiegare, oltre alla teoria, anche gli strumenti pratici con i quali hanno provato a venirne a capo. E i risultati ottenuti con le previsioni sulle primarie nel Pd o con le percentuali del M5S alle Politiche del 2013 parlano a loro favore. Nessuno di noi a Wired li prese sul serio quando gli autori ci comunicarono che secondo loro il movimento sarebbe arrivato oltre il 20%, mentre gli altri istituti di ricerca davano percentuali molto più basse. Come sia andata lo sappiamo.

Infine, un ulteriore motivo di attenzione per questo volume è che vi si racconta di un'area di studi innovativi oggi trasformata in un'impresa. Voices from the blogs è una Srl, spin off dell'Università di Milano. Raro che l'accademia diventi azienda, in Italia.

Ci saranno sempre eventi che sfuggono alla prevedibilità, come immaginava Asimov con la sua psicostoria e come riconoscono gli autori, consapevoli dei limiti odierni di questa disciplina. "Don't panic": non facciamoci prendere dal panico, consigliava Adams agli autostoppisti galattici. Ci saranno nuove tecniche da sperimentare per affinare i risultati delle previsioni, nuove aziende e professioni da inventare per cogliere altrettante opportunità, nuovi libri da scrivere. Se saranno come questo, sarà piacevole leggerli.

Massimo Russo
Direttore di Wired Italia

Lo scorso febbraio stavo attraversando, in treno, le pianure innevate dell'Emilia-Romagna e della Lombardia, per partecipare ad una conferenza sui social media che si sarebbe svolta a Milano. Solo un giorno prima, al Campus della New York University a Firenze, avevo partecipato ad una serie di dibattiti sulle elezioni presidenziali statunitensi in cui esprimevo il mio pensiero rispetto al ruolo crescente dei social media nella politica americana.

Il mio intervento, almeno così mi auguro, ha fornito interessanti spunti di riflessione, ma l'aspetto che ha suscitato in me il maggior interesse è stato sicuramente il lavoro di Voices from the Blogs (VfB). Questo gruppo di ricerca ha costruito e sviluppato un metodo di sentiment analysis delle informazioni postate sulla Rete, in grado di prevedere in modo sorprendente, tra l'altro, i risultati delle elezioni.

In quegli stessi giorni la campagna elettorale italiana si avviava alla sua conclusione, e le leggi italiane vietavano ai fondatori di VfB la pubblicazione delle loro stime di voto. Ho avuto modo di vedere i loro dati e di confrontarli, ad urne chiuse, con i voti reali. Il processo elettorale italiano è a dir poco labirintico, ma nonostante questo le previsioni di VfB sono state straordinariamente vicine ai risultati delle elezioni per la Camera dei Deputati.

Personalmente, in quanto americano, ero ancora più stupito dalla performance di VfB nel prevedere l'esito della sfida Obama-Romney. La notte prima del voto Usa, il blog di VfB ha pubblicato le stime realizzate attraverso questa metodologia di previsione basata su Twitter, pronosticando una vittoria del presidente uscente con un margine del 3,5% nel voto popolare. Il distacco effettivo fu di 3,9 punti. La media dei sondaggi realizzata da Real Clear Politics invece prevedeva una sottilissima vittoria di Obama, con un vantaggio dello 0,7%. Se Mitt Romney avesse prestato attenzione a Voices from the Blogs non sarebbe rimasto così sconcertato dall'esito del voto.

Questa tecnica di analisi è solo una delle ultime rivoluzioni nel rapporto tra politica e social media, che si è avviato negli Stati Uniti e da lì ha contagiato altri paesi tanto da permettere a movimenti di protesta come quello promosso da Beppe Grillo di rafforzarsi fino ad ottenere il 25% dei voti. Ironia della sorte, la rivoluzione del web non ha portato alcun beneficio ad Al Gore che, a dispetto della caricatura che ne fanno i media, introdusse la legislazione che avviò la commercializzazione della tecnologia DARPA al Dipartimento della Difesa Usa, tecnologia che ci ha consentito di avere accesso ad Internet.

Nel 2004 i social media hanno permesso a Howard Dean di raccogliere 40 milioni di dollari già prima del caucus dell'Iowa, che rappresentava il primo grande evento nella corsa delle primarie. Successivamente, lo staff che organizzò la campagna presidenziale di Kerry trasse ispirazione da questo strumento e infatti l'utilizzo dei social media giocò un ruolo cruciale nella raccolta fondi, permettendo a John Kerry di incassare circa 250 milioni di dollari prima ancora che la convention del Partito Democratico ufficializzasse la sua nomination, il che è stupefacente, perché garantì a Kerry la possibilità di essere competitivo con Bush dal punto di vista dei finanziamenti raccolti, almeno fino al momento della convention. Alla luce di questo, Kerry stesso ha ammesso l'errore, forse decisivo, di non aver fatto ricorso ai social media per finanziare anche la fase finale della campagna elettorale, rinunciando di fatto alla possibilità di raccogliere centinaia di milioni di dollari attraverso i canali social.

Obama non ripeté lo stesso errore nel 2008, e fu in grado, proprio grazie a questo, di battere John McCain in termini di raccolta fondi, sopravanzandolo di quasi 300 milioni di dollari. Ma nel 2008 l'utilizzo dei social media non si limitò alla funzione di raccolta fondi. Al contrario, i social media divennero il canale per una comunicazione a "doppio senso", sia bottom-up che top-down, tra lo staff della campagna e gli elettori. Nel 2012 i canali di comunicazione "non-tradizionali" hanno permesso a Obama di effettuare il cosiddetto "micro-targeting", sviluppando la più imponente ed efficiente operazione mai condotta, per convincere e mobilitare l'elettorato prima del voto. La combinazione tra campagna sul territorio e sui social media ha poi fatto la differenza. Romney, al contrario, scelse di contattare il proprio elettorato attraverso un network (chiamato "the Whale") basato sulla telefonia cellulare, network che però è collassato proprio nel momento più importante: nel giorno delle elezioni.

Questo libro affronta tutta un'altra serie di promettenti sviluppi nel campo dei social media. Le analisi di VfB, effettuate via Twitter, sono in grado di predire il risultato di un'elezione. Ma si tratta solo di questo? O invece questo strumento può essere utilizzato per effettuare una sorta di "diagnostica" delle strategie elettorali, permettendo così di modificare, smussare e indirizzare i messaggi della campagna,

alterando in questo modo l'esito della sfida? Per il futuro mi aspetto esattamente che accada questo, e molto altro.

Esistono già applicazioni commerciali che utilizzano tecniche simili a quelle impiegate dal gruppo di Voices from the Blogs (ma non sono sicuro, in realtà, che ci sia bisogno di VfB per scoprire che i passeggeri hanno un'immagine più negativa, ad esempio, di Alitalia dopo aver utilizzato il servizio).

Riguardo alle campagne elettorali invece è tutta un'altra storia. Gli analisti e gli staff della campagna hanno una sete inesauribile di nuovi metodi utili a studiare, raggiungere, persuadere e mobilitare l'elettorato. La tecnologia non rimpiazzerà mai il messaggio, ma può indirizzarlo, amplificarlo, personalizzarlo. Le tecniche che VfB ha sviluppato sarà utilizzata ed applicata in modi nuovi e diversi, e presto o tardi permetterà ai candidati di scoprire non tanto se sono in testa nelle intenzioni di voto, ma piuttosto cosa fare per riuscire a vincere.

Robert Shrum
Senior Fellow e Clinical Professor
presso la Robert F. Wagner School of Public Service della New York University

Prefazione

*“Quarantadue!” urlò Loonquawl.
“Questo è tutto ciò che sai dire dopo un lavoro di sette milioni e mezzo di anni?”
“Ho controllato molto approfonditamente” disse il computer, “e questa è
sicuramente la risposta.*

*Ad essere sinceri, penso che il problema sia che voi non abbiate mai saputo
veramente qual è la domanda”*

Douglas Noël Adams, Guida galattica per gli autostoppisti

La grande diffusione dei social media ed il loro ruolo nelle società contemporanee, rappresentano una delle novità più interessanti di questi ultimi anni, tanto da aver catturato l'interesse di ricercatori, giornalisti, imprese, movimenti e governi. La crescita della facilità di accesso alle informazioni che la rete permette e l'opportunità potenziale di comunicare con una vasta platea ad un costo praticamente nullo viene interpretata spesso come un passo verso una democratizzazione del discorso pubblico.¹ Questo, ovviamente, non significa trascurare gli squilibri intrinseci che i nuovi media portano con sé – come evidenziato ad esempio dalla distribuzione a legge di potenza (o distribuzione *power-law*) della rete, in cui ad un numero contenuto di siti o utenti che ricevono moltissime visite, fa da contraltare una vasta maggioranza di altri siti o utenti con un numero di lettori estremamente limitato.² Detto questo, però, la densa interconnessione che spesso si crea tra chi è attivo in rete genera uno spazio di discussione che è in grado di motivare e coinvolgere gli individui in una più ampia agorà, collegando tra di loro persone con obiettivi comuni e facilitando così forme disperate di azione collettiva.³ Questo, a sua volta, dà vita a ciò che viene definito “individualismo in rete”⁴: invece di contare sempre su una singola comunità di riferimento, grazie ad internet e, soprattutto, ai social media, diventa oggi possibile muoversi tra più persone e risorse, spesso eterogenee tra di loro, a seconda

¹ Farrell H (2012) The Internet's Consequences for Politics. *Annual Review of Political Science* 15:35–52.

² Adamic LA, Huberman BA (2000) Power-Law Distribution of the World Wide Web. *Science* 287(5461):2115. URL: <http://www.sciencemag.org/content/287/5461/2115.full>.

³ Valenzuela S, Park N, Kee KF (2009) Is there social capital in a social network site? Facebook use, and college students' life satisfaction, trust, and participation. *Journal of Computer-Mediated Communication* 14(4):875–901.

⁴ Wellman B (2001) Physical Place and Cyber Place: The Rise of Networked Individualism. *International Journal of Urban and Regional Research* 25(2):227–252.

delle situazioni via via sostenute dal singolo utente, selezionando quelle più adatte a risolvere particolari esigenze o ad approfondire determinati interessi.

Pur con tutte le cautele necessarie, sembra dunque che i social media stiano dando vita ad una rivoluzione digitale. E l'aspetto più interessante di questo cambiamento non è tanto (o unicamente) legato alle possibilità di favorire partecipazione politica e attivismo, come da più parti viene sostenuto. La vera *social revolution* è quella che investe le vite di ogni singolo individuo. È la libertà di esprimersi, di avere uno spazio proprio in cui essere sé stessi, o ciò che si vorrebbe essere, con pochi limiti e barriere. La rivoluzione dei social media è allora quella di poter raccontare le proprie emozioni ed opinioni non solo a sé stessi, quanto, e soprattutto, a chi ci circonda, interagendo con loro, aprendo reciprocamente una finestra sui rispettivi mondi, curiosando con occhio più o meno indiscreto sulle *vite degli altri*.

E tutto questo accade, paradossalmente, mentre viviamo in una società in cui si fatica sempre più a conoscere i nomi dei vicini di casa, e in cui il diritto alla privacy diventa un imperativo a cui sottostare, salvo poi comunicare a tutto il mondo (rigorosamente on-line) qualunque cosa: amori, delitti, giorni indimenticabili e fallimenti quotidiani. Perché in effetti sui social media (o, meglio, su quei social media che sono anche social network, come vedremo) si finisce per raccontare tutta (o quasi) la propria vita: dalla felicità per la nascita di un figlio, alla rabbia per un treno in ritardo, dallo shopping pre-natalizio alla scelta di voto fatta nel *segreto* della cabina elettorale.

Non c'è allora da stupirsi se da più parti si sia incominciato a discutere delle modalità attraverso cui utilizzare al meglio questo *mare magnum* di informazioni. Perché i dati presenti in rete, se opportunamente raccolti e analizzati, permettono non solo di capire e spiegare molti fenomeni sociali complessi, ma anche, e persino, di prevederli. La previsione, sia quella fatta in tempo reale che quella relativa ad eventi futuri, è in effetti una delle frontiere più seducenti del mondo *social*.

E così, secondo alcuni, acquista senso l'idea di concepire i social media come una pluralità di "antenne" interconnesse tra di loro, che agiscono come una sorta di "cervello collettivo" in grado di cogliere, a volte meglio di altre alternative, i trend che si dipanano continuamente intorno a noi. Una proprietà emergente di un vero e proprio sistema complesso,⁵ che nasce dall'aggregazione di tutte le sue componenti, invece che essere riconducibile ad una sola singola parte.

Nel Cap. 1 partiremo proprio da questo aspetto. Dopo aver discusso brevemente le caratteristiche e dato qualche cifra in merito alla diffusione dei social media in Italia e nel mondo, riassumeremo alcune aree di analisi per concentrarci poi su una rassegna dei più interessanti esempi di previsione fatti utilizzando i social media, spaziando dall'economia alla epidemiologia, dal marketing alla politica.

Nel Cap. 2 effettueremo il passo successivo, discutendo più in dettaglio delle diverse tecniche utilizzate finora per analizzare la rete, presentando e confrontando in tal senso i diversi metodi per fare *sentiment analysis*, ovvero per monitorare l'umore di chi scrive sui social media rispetto ai più svariati argomenti, cercando di estrarne un significato operativo. Pur nella sua brevità, questo capitolo fornirà una

⁵ Axelrod R (1997) *The Complexity of Cooperation*. Princeton University Press, Princeton; Curini L (2009) *Gli agent-based models: come modellare la complessità*. Quaderni di Scienza Politica 3:517-531.

bussola ragionata attraverso cui orientarsi in una letteratura che di anno in anno sta crescendo in modo esponenziale. Questo ci porterà poi a presentare i vantaggi di un nuovo metodo, la *integrated Sentiment Analysis (iSA)*, basata su una ottimizzazione dell'algoritmo introdotto originariamente da Hopkins e King.⁶

Nel Cap. 3 attraverso questa tecnica analizzeremo i social media per misurare un'emozione che risulta quanto mai sfuggente, ma di cui si è fatto un gran discutere in questi ultimi anni come possibile alternativa (o complemento) alle misure più strettamente economiche del benessere di una collettività, ossia la felicità. In particolare, focalizzandoci su oltre 40 milioni di messaggi postati su Twitter, studieremo la felicità degli italiani nel corso di tutto il 2012, mostrandone l'andamento, giorno per giorno, nelle 110 province. A partire da questi dati svilupperemo poi una analisi econometrica per provare a spiegare quali fattori facciano crescere o diminuire la felicità in Italia.

Nel Cap. 4 ci concentreremo invece sulla possibilità di prevedere un evento concreto, ovvero i risultati delle elezioni politiche, attraverso ciò che viene pubblicato sui social media. Per farlo analizzeremo una pluralità di tornate elettorali che hanno avuto luogo tra il 2012 e il 2013 in contesti molto diversi tra loro: dalle elezioni presidenziali e legislative francesi alle presidenziali Usa, dalle primarie del centrosinistra, fino alle elezioni politiche italiane. In ciascun caso confronteremo le nostre previsioni con i dati dei sondaggi e con l'esito dell'urna. I risultati, come vedremo, sono decisamente promettenti sia per quanto riguarda la capacità dei social media di narrare fedelmente ed in tempo reale l'evoluzione ed i trend di una campagna elettorale, sia in termini di accuratezza rispetto all'anticipazione del risultato finale.

Nelle conclusioni, infine, tratteremo nuovi indirizzi di ricerca oltre ad alcuni possibili sviluppi e applicazioni rese possibili dall'analisi della rete. Prenderemo prima in considerazione il ruolo di social media come strumento che può permettere, almeno in potenza, di migliorare l'*accountability* e la *responsiveness* dei governanti verso i governati, per poi tornare al discorso da cui siamo partiti, ovvero all'idea dei social media come strumento di *e-Democracy*.

Terminiamo questa introduzione con un doveroso ringraziamento a tutti coloro che hanno contribuito alla realizzazione di questo lavoro. Primo tra tutti Giuseppe Porro, che ha contribuito ad avviare due anni fa, assieme agli autori del presente volume, *Voices from the Blogs* (VfB: <http://voicesfromtheblogs.com>), un progetto di ricerca dell'Università degli Studi di Milano volto a monitorare sistematicamente i commenti pubblicati sulla rete e sui social media, e le cui analisi e metodologie sono alle fondamenta delle pagine qui riportate. Ringraziamo poi Irina Iasinovschi, che ha dato un contributo rilevante per la raccolta dei dati pubblicati nel Cap. 1 di questo lavoro, nonché per la parte grafica, e Luciano Canova per la discussione ed i preziosi suggerimenti in relazione al tema della felicità.

Un altro doveroso ringraziamento va a Wired Italia con cui abbiamo lungamente collaborato in questi mesi ed in particolare a Federico Ferrazza che ha sostenuto e diffuso ad un vasto pubblico molte delle analisi che sono riportate in questo libro

⁶ Hopkins D, King G (2010) A Method of Automated Nonparametric Content Analysis for Social Science. *American Journal of Political Science* 54(1):229–247.

(dedicarci la storia di copertina del numero di Wired di gennaio 2012 è stato un gesto quasi “avventato”... ma molto apprezzato). Ma non saremmo mai arrivati a questo punto senza il sostegno di Renato Mattioni, segretario della Camera di Commercio di Monza-Brianza, con il quale abbiamo pubblicato un primo libro in cui sono stati analizzati i tweet dei milanesi, toccando argomenti che vanno dalla politica, al design, al lavoro.⁷ Un grande ringraziamento va anche ad Emanuela Croci, Mario Barone, Massimo Donelli, Andrea Vicari, e al Consolato Generale degli Stati Uniti d’America di Milano, in particolare a Francesca Bettelli e Donatello Osti, assieme a cui abbiamo organizzato diverse iniziative.

Molte delle analisi qui riportate sono state pubblicate sul sito del Corriere della Sera, che ringraziamo perché da oltre un anno ospita il nostro blog (<http://sentimeter.corriere.it>), e per aver pubblicato all’interno dei suoi “speciali” i dati relativi alle elezioni presidenziali americane, alle elezioni primarie del centro-sinistra e alle elezioni politiche (almeno fino allo stop imposto da AgCom...). In particolare, ringraziamo Paolo Ottolina, Paolo Rastelli e Giovanni Angeli. Ringraziamo inoltre Luca Tremolada del Sole24Ore con il quale abbiamo lavorato su diversi temi tecnologici come il lancio del nuovo iPad e del nuovo iPhone.

Altri ringraziamenti assai sentiti vanno a Simone Spetia di Radio 24 che ci ha più volte ospitato per parlare di alcuni degli argomenti discussi in questo libro, ma non solo, all’interno della trasmissione Votantonio. Un ringraziamento è dovuto anche a Walter Galbiati e a Economia e Finanza di Repubblica.it, così come al team di Radio 2 – Miracolo Italiano, con cui abbiamo passato l’estate del 2013 a parlare di felicità.

È doveroso ringraziare naturalmente anche l’Università degli Studi di Milano, e in particolare la squadra di UNIMITT (Innovazione e Trasferimento Tecnologico Università degli Studi di Milano) nelle persone di Chiara del Balio e Roberto Tiezzi che hanno supportato la nascita dello spin-off Voices from the Blogs srl, così come il Rettore Gianluca Vago, la cui firma ha dato ufficialmente il via all’iniziativa, in una data molto particolare: 12/12/12. Allo stesso modo vogliamo ringraziare Marco Giuliani e Franco Donzelli, rispettivamente direttori del Dipartimento di Scienze Sociali e Politiche e del Dipartimento di Economia, Management e Metodi Quantitativi dell’Università degli Studi di Milano, per aver sostenuto lo sviluppo di VfB collaborando anche all’organizzazione di diversi seminari tematici. Un ringraziamento va anche a diversi colleghi e assegnisti di ricerca che sono attualmente afferenti ai suddetti dipartimenti tra cui Mauro Barisione e Marco Mainenti, o lo sono stati in passato, come Vincenzo Memoli.

Un particolare ringraziamento va all’amico e collega Gary King, dell’Università di Harvard, per averci introdotto nel lontano 2010 all’analisi dei social media, nonché all’istituto ISPO nelle persone di Renato e Ludovico Manheimer per la disponibilità accordata ad effettuare ricerche congiunte e analisi di validazione dei nostri risultati.

Questo libro non sarebbe poi stato possibile senza l’aiuto, il sostegno e l’entusiasmo di molti studenti, laureandi e laureati della Facoltà di Scienze Politiche, Eco-

⁷ Ceron A, Curini L, Iacus S, Mattioni R, Porro G (2012) #Milano-Brianza in un tweet: lavoro, politica, partecipazione. Guerini e Associati, Milano.

nomiche e Sociali dell'Università degli Studi di Milano. Vogliamo qua ricordare, in stretto ordine alfabetico: Vito Andreana, Agnese Barni, Giulia Baronio, Cinzia Besana, Alice Blangero, Angelo Boccato, Filippo Caracciolo, Luca Castelli, Matteo Ceccarelli, Barbara Colombini, Francesca Condoleo, Alessandra Caterina Cremonesi, Alessandro Del Tredici, Stefano Doronzo, Gianluca Gaiga, Alberto Galbusera, Ilaria Locorotondo, Marco Moggia, Luca Noris, Giovanni Paini, Benedetta Pinò, Francesco Russo, Salvatore Salamone, Eliza Ungaro.

Gli ultimi e più sentiti ringraziamenti vanno a Spoletta, la nostra mascotte, per aver saputo illuminare il gruppo di VfB sin dalla sua fondazione e alle persone care che ci hanno supportato e continuano a farlo nonostante le assenze, le notti insonni e le alzatacce passate ad analizzare i social media.

Milano, settembre 2013

Andrea Ceron
Luigi Curini
Stefano M. Iacus

1	Perché studiare i social media	1
1.1	I social media: caratteristiche e definizioni	1
1.2	“Utenti della Rete, unitevi!”: alcuni numeri sulla diffusione dei social media	4
1.3	Principali direzioni di ricerca sui social media	9
1.3.1	Social media e previsioni	12
1.4	I vantaggi dell’analisi via Twitter	17
	Riferimenti web	18
	Riferimenti bibliografici	22
2	Opinion Mining e integrated Sentiment Analysis (iSA)	27
2.1	Dall’analisi del linguaggio alle opinioni	27
2.1.1	Analisi quantitativa e analisi qualitativa dei testi	27
2.1.2	I principi fondamentali dell’analisi testuale	28
2.2	L’analisi dei testi in pratica	31
2.2.1	Come rendere il testo digeribile ad un modello statistico: lo <i>stemming</i>	31
2.2.2	Le famiglie di tecniche di analisi testuale: lo <i>scoring</i>	34
2.2.3	Pregi e difetti del <i>tagging</i> automatico e umano	35
2.2.4	Metodi di classificazione testuale	36
2.2.5	Tecniche di <i>clustering</i>	36
2.2.6	Topic models	37
2.2.7	Classificazione individuale e aggregata: il contributo di Hopkins e King	38
2.2.8	Perché si riduce l’errore?	41
2.2.9	Il problema del rumore	41
2.2.10	Quanto deve essere grande il <i>training set</i> nel metodo <i>iSA</i> ?	42
2.2.11	Segnale forte, segnale debole e stime vincolate	42
2.2.12	I vantaggi di <i>iSA</i>	43
2.2.13	Integrazione di metodi <i>supervised</i> e tecniche di <i>scoring</i>	44
2.2.14	Altri approcci all’analisi dei testi	44