# Social Media: Friend or Foe of Natural Language Processing?

## Tim Baldwin



THE UNIVERSITY OF
MELBOURNE

# Talk Outline

# What is Social Media?

- According to Wiktionary (21/8/2012), social media is:

    *Interactive forms of media that allow users to interact with and publish to each other, generally by means of the Internet.*

- While social media sites have strong support for multimedia content, text is still very much a core data type

# Social Media Include ...

💡 **Social Networking sites**

posts, friends/circles, "likes", shares, events, photos, comments, geotags, ...



**2010** 500 million users

**Source(s):** http://mashable.com/2011/02/04/facebook-7th-birthday/

# Social Media Include ...

> 💡 **Micro-blogs**
> posts, followers/followees, shares, hashtagging, geotags, ...

# Social Media Include ...

## Web user forums

posts, threading, followers/followees, ...



**Source(s):**

http://http://forums.cnet.com/7723-6617_102-570394/ubuntu-running-minecraft/

# Social Media Include ...

💡 **Wikis**
| posts, versioning, linking, tagging, ...



**Source(s):** `http://en.wikipedia.org/wiki/Social_media`

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- Edited text

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

**? How different?**

Bigram LM Perplexity:

|  | BNC→ | Twitter$_1$→ | Twitter$_2$→ |
|---|---|---|---|
| →BNC | 185 | 1553 | 1528 |
| →Twitter$_1$ | 4082 | 260 | 887 |
| →Twitter$_2$ | 4953 | 938 | 274 |

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- Static data

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

**? Challenges of Streaming Data**
require throughput guarantees
batch vs. streamed processing of data (e.g. for topic modelling)
potential need for "incremental" models

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- Long(ish) documents; plenty of context

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

> **? Document Context**
> Hard to adjust document-level priors when little context

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*
- All context is language context

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*
- *Little language, potentially lots of other context*

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

**? Priors, priors everywhere**
user priors
user-declared metadata priors
location priors
social network-based priors
hashtag priors
timezone priors
implicit social networks (retweets, user mentions, ...)
⋮

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*
- *Little language, potentially lots of other context*
- Well-defined domain/genre

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*
- *Little language, potentially lots of other context*
- *All over the place*

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*
- *Little language, potentially lots of other context*
- *All over the place*
- Sentence tokenisation

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*
- *Little language, potentially lots of other context*
- *All over the place*
- *What's a sentence?*

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*
- *Little language, potentially lots of other context*
- *All over the place*
- *What's a sentence?*
- Grammaticality

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*
- *Little language, potentially lots of other context*
- *All over the place*
- *What's a sentence?*
- *Yer what?*

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*
- *Little language, potentially lots of other context*
- *All over the place*
- *What's a sentence?*
- *Yer what?*
- Most of what glitters is English (and if your method can handle one language, it can handle 'em all)

# Properties of Social Media Data

(NLP "ideal" → *actuality*)

- *Unedited text*
- *Streamed data*
- *Short documents; v. little context*
- *Little language, potentially lots of other context*
- *All over the place*
- *What's a sentence?*
- *Yer what?*
- *Anything goes — lots of languages, multilingual documents, ad hoc spelling, mix of language and markup ... language anarchy!*

# Observation/Questions

- Much of the work that is currently being carried out over social media data doesn't make use of NLP

# Observation/Questions

- Much of the work that is currently being carried out over social media data doesn't make use of NLP
    - Are NLP methods not suited to social media analysis?

# Observation/Questions

- Much of the work that is currently being carried out over social media data doesn't make use of NLP
    - Are NLP methods not suited to social media analysis?
    - Is social media data too challenging for modern-day NLP?

# Observation/Questions

- Much of the work that is currently being carried out over social media data doesn't make use of NLP
  - Are NLP methods not suited to social media analysis?
  - Is social media data too challenging for modern-day NLP?
  - Are simple term search-based methods sufficient for social media analysis, i.e. is NLP *overkill* for social media?

# Observation/Questions

- Much of the work that is currently being carried out over social media data doesn't make use of NLP
  - Are NLP methods not suited to social media analysis?
  - Is social media data too challenging for modern-day NLP?
  - Are simple term search-based methods sufficient for social media analysis, i.e. is NLP *overkill* for social media?

- Is social media data is the friend or foe of NLP?

# Possible Ways Forward

# Possible Ways Forward

- "Adapt" the data to the NLP tools through preprocessing of various forms

# Possible Ways Forward

- "Adapt" the data to the NLP tools through preprocessing of various forms
- "Adapt" the NLP tools to the data through "domain" (de-)adaptation

# Talk Outline

# Preprocessing

- Basic premise: the cleaner/richer the data, the easier it is to process/better quality the predictions that arise from it
- Overarching constraint: any preprocessing has to be able to keep pace with the torrent of streamed data ... although many of the models we use can be learned off-line

# Language Identification: Task

- Language identification (langid) = prediction of the language(s) a given message is authored in

> **? Example**
> *karena ada rencana ke javanet, maka siapkan link dolodan, di bookmark, ready to be a bandwidth killer.. siap siaplah javanet, im coming..*
> Language(s): ?

# Language Identification: Task

- Language identification (langid) = prediction of the language(s) a given message is authored in

> **? Example**
>
> *karena ada rencana ke javanet, maka siap-kan link dolodan, di bookmark, ready to be a bandwidth killer.. siap siaplah javanet, im coming..*
>
> Language(s): MS,EN

# Language Identification: Method

- Outline of the basic approach:

**Source(s):** Baldwin and Lui [2010], Lui and Baldwin [2011, 2012]

# Language Identification: Method

- Outline of the basic approach:
  1. represent each document as a set of byte $n$-grams of varying $n$

**Source(s):** Baldwin and Lui [2010], Lui and Baldwin [2011, 2012]

# Language Identification: Method

- Outline of the basic approach:
  1. represent each document as a set of byte *n*-grams of varying *n*
  2. across a range of datasets, identify *n*-grams that are correlated with language and *not* dataset

**?LD**

$$\mathcal{LD}^{all}(t) = \mathcal{IG}_{lang}^{all}(t) - \mathcal{IG}_{domain}(t)$$

# Language Identification: Method

- Outline of the basic approach:
  1. represent each document as a set of byte *n*-grams of varying *n*
  2. across a range of datasets, identify *n*-grams that are correlated with language and *not* dataset
  3. learn log likelihoods for each term and class from training data: $logP(t_j|c_i)$

**Source(s):** Baldwin and Lui [2010], Lui and Baldwin [2011, 2012]

# Language Identification: Method

- Outline of the basic approach:
  1. represent each document as a set of byte *n*-grams of varying *n*
  2. across a range of datasets, identify *n*-grams that are correlated with language and *not* dataset
  3. learn log likelihoods for each term and class from training data: $logP(t_j|c_i)$
  4. classify a test document using multinomial naive Bayes over the $\mathcal{LD}$ features

**Source(s):** Baldwin and Lui [2010], Lui and Baldwin [2011, 2012]

# Language Identification: Accuracy

- Comparative evaluation over pre-existing Twitter LangID datasets:

|      | langid.py | | LangDetect | | CLD | |
|------|-----------|---------|-------|----------|------|----------|
|      | Accuracy | docs/s | ΔAcc | Slowdown | ΔAcc | Slowdown |
| T-BE | 0.941 | 367.9 | −0.016 | 4.4× | −0.081 | 0.7× |
| T-SC | 0.886 | 298.2 | −0.038 | 2.9× | −0.120 | 0.2× |

- Impact on bigram LM Perplexity:

|      | BNC→ | | Twitter-EN$_1$→ | | Twitter-EN$_2$→ | |
|------|------|--|------|--|------|--|
| →BNC | 185 | | 1170 | (−383) | 1108 | (−420) |
| →Twitter-EN$_1$ | 1528 | (−2554) | 215 | | 416 | (−471) |
| →Twitter-EN$_2$ | 1620 | (−333) | 469 | (−469) | 228 | |

**Source(s):** Baldwin and Lui [2010], Lui and Baldwin [2011, 2012]

# Language Identification: Research Challenges

- We are very good at monolingual language identification, but what about multilingual documents?

# Language Identification: Research Challenges

- We are very good at monolingual language identification, but what about multilingual documents?
  - multi-label language identification (*what language(s) is a document in*)

# Language Identification: Research Challenges

- We are very good at monolingual language identification, but what about multilingual documents?
  - multi-label language identification (*what language(s) is a document in*)
  - language segmentation (*which parts of what messages correspond to what languages?*)

# Language Identification: Research Challenges

- We are very good at monolingual language identification, but what about multilingual documents?
  - multi-label language identification (*what language(s) is a document in*)
  - language segmentation (*which parts of what messages correspond to what languages?*)
- How can we determine when we aren't sure/don't recognise the language(s)?

# Lexical Normalisation: Task

- Lexical normalisation = "spell-correct" (English) messages to "canonical" lexical form:

> **? Example**
>
> *If you a GIrl and you dont kno how to Cook yo bf should Leave you rite away*
>
> ↓
>
> *If you a girl and you don't know how to cook your boyfriend should leave you rite away*

**Source(s):** Han and Baldwin [2011], Han et al. [2012], Gouws et al. [2011], Liu et al. [2011, 2012]

# Lexical Normalisation: Method

- Outline of approach:

# Lexical Normalisation: Method

- Outline of approach:
    1. pre-learn (OOV,IV) word pairs from microblog data

**Source(s):** Han and Baldwin [2011], Han et al. [2012]

# Lexical Normalisation: Method

- Outline of approach:
  1. pre-learn (OOV,IV) word pairs from microblog data
  2. lexical normalisation by simple dictionary lookup

**Source(s):** Han and Baldwin [2011], Han et al. [2012]

# Lexical Normalisation: Method

- Outline of approach:
  1. pre-learn (OOV,IV) word pairs from microblog data
  2. lexical normalisation by simple dictionary lookup
- Learning the normalisation dictionary:

**Source(s):** Han and Baldwin [2011], Han et al. [2012]

# Lexical Normalisation: Method

- Outline of approach:
  1. pre-learn (OOV,IV) word pairs from microblog data
  2. lexical normalisation by simple dictionary lookup
- Learning the normalisation dictionary:
  1. Extract (OOV, IV) pairs based on distributional similarity.

**Source(s):** Han and Baldwin [2011], Han et al. [2012]

# Lexical Normalisation: Method

- Outline of approach:
  1. pre-learn (OOV,IV) word pairs from microblog data
  2. lexical normalisation by simple dictionary lookup
- Learning the normalisation dictionary:
  1. Extract (OOV, IV) pairs based on distributional similarity.
  2. Re-rank the extracted pairs by string similarity.

**Source(s):** Han and Baldwin [2011], Han et al. [2012]

# Lexical Normalisation: Method

- Outline of approach:
  1. pre-learn (OOV,IV) word pairs from microblog data
  2. lexical normalisation by simple dictionary lookup
- Learning the normalisation dictionary:
  1. Extract (OOV, IV) pairs based on distributional similarity.
  2. Re-rank the extracted pairs by string similarity.
  3. Select the top-$n$ pairs for inclusion in the normalisation lexicon.

**Source(s):** Han and Baldwin [2011], Han et al. [2012]

# Lexical Normalisation: Results

- Lexical normalisation results:

| Method | Precision | Recall | F-Score |
|---|---|---|---|
| S-dict | 0.700 | 0.179 | 0.285 |
| HB-dict | 0.915 | 0.435 | 0.590 |
| GHM-dict | **0.982** | 0.319 | 0.482 |
| HB-dict+GHM-dict+S-dict | 0.847 | 0.630 | **0.723** |

Ultimately: dictionary combination works best

**Source(s):** Han and Baldwin [2011], Han et al. [2012]

# Lexical Normalisation: Results

- Lexical normalisation results:

| Method | Precision | Recall | F-Score |
|---|---|---|---|
| S-dict | 0.700 | 0.179 | 0.285 |
| HB-dict | 0.915 | 0.435 | 0.590 |
| GHM-dict | **0.982** | 0.319 | 0.482 |
| HB-dict+GHM-dict+S-dict | 0.847 | 0.630 | **0.723** |

Ultimately: dictionary combination works best

- Impact on POS tagging:

| Tagger | Text | % accuracy | # correct tags |
|---|---|---|---|
| $POS_{Stanford}$ | original | 68.4 | 4753 |
| $POS_{Stanford}$ | normalised | 70.0 | 4861 |
| $POS_{twitter}$ | original | 95.2 | 6819 |
| $POS_{twitter}$ | normalised | 94.7 | 6780 |

**Source(s):** Han and Baldwin [2011]; Han et al. [2012]

# Other Instances of Preprocessing

- User/message geolocation
- Identification of "high-utility" messages
- Social media user profiling
- Credibility analysis

# Talk Outline

1. Social Media and Natural Language Processing

2. Bringing the Data to NLP

3. **Bringing NLP to the Data**

4. Concluding Remarks

# Instances of Social Media-adapted NLP tools

- CMU Twitter POS tagger: Twitter-tuned, coarse-grained POS tagset
- Self-training parser adaptation for social media data
- Named Entity Recognition for Twitter

**Source(s):** Gimpel et al. [2011], Foster et al. [2011], Ritter et al. [2011]

# The Grand Challenge

- Social media data is highly temporal in nature, and models constantly need updating/de-adaptation
- Often in social media analysis, people are interested in finding the *unknown* (e.g. novel event types, new products)

# Talk Outline

# Friend or Foe?

# Friend or Foe?

- If as NLPers we cherish a challenge, there is no question that social media is our friend

# Friend or Foe?

- If as NLPers we cherish a challenge, there is no question that social media is our friend
- If we simplistically apply models trained on "traditional" datasets to social media, it is very much a foe … and evermore shall be so!

# Friend or Foe?

- If as NLPers we cherish a challenge, there is no question that social media is our friend
- If we simplistically apply models trained on "traditional" datasets to social media, it is very much a foe … and evermore shall be so!
- Social media also opens up immediate opportunities in terms of integrated multimodal analysis (links to image, video content); if we can harness this content, social media is again our friend (more context/better disambiguation)

# NLP and Social Media

- Is NLP overkill for social media analysis?

**Source(s):** Ritterman et al. [2009], Sakaki et al. [2010]

# NLP and Social Media

- Is NLP overkill for social media analysis?
- Much of the work on social media analysis is based on analysis of a pre-defined trend (e.g. election outcome prediction, flu outbreak tracking, earthquake detection)

**Source(s):** Ritterman et al. [2009], Sakaki et al. [2010]

# NLP and Social Media

- Is NLP overkill for social media analysis?
- Much of the work on social media analysis is based on analysis of a pre-defined trend (e.g. election outcome prediction, flu outbreak tracking, earthquake detection)

    … and perhaps NLP is overkill

**Source(s):** Ritterman et al. [2009], Sakaki et al. [2010]

# NLP and Social Media

- Is NLP overkill for social media analysis?
- Much of the work on social media analysis is based on analysis of a pre-defined trend (e.g. election outcome prediction, flu outbreak tracking, earthquake detection)

    ... and perhaps NLP is overkill

- That is not to say there aren't a myriad of applications which can't be described with simple keywords for which NLP is vital (e.g. novel event detection, disaster management)

**Source(s):** Ritterman et al. [2009], Sakaki et al. [2010]

# NLP and Social Media

- Is NLP overkill for social media analysis?
- Much of the work on social media analysis is based on analysis of a pre-defined trend (e.g. election outcome prediction, flu outbreak tracking, earthquake detection)

  ... and perhaps NLP is overkill

- That is not to say there aren't a myriad of applications which can't be described with simple keywords for which NLP is vital (e.g. novel event detection, disaster management)

  ... and perhaps the bottleneck is instead NLP accessibility

**Source(s):** Ritterman et al. [2009], Sakaki et al. [2010]

# Final Words

- Social media is hip ... but also big and hairy, and poses both challenges and opportunities for NLP
- Ongoing work on a myriad of technologies/tasks relating to social media analysis, progressively making social media more "NLP accessible"
- There is plenty to be done ... come and join us!

# Taking Credit for a Cast of Thousands

- This is joint work with Paul Cook, Bo Han, Aaron Harwood, Shanika Karunasekera, Su Nam Kim, Marco Lui, David Martinez, Joakim Nivre, Richard Penman, Li Wang, ...

# References I

Timothy Baldwin and Marco Lui. Language identification: The long and the short of the matter. In *Proc. of Human Language Technologies: The 11th Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL HLT 2010)*, pages 229–237, Los Angeles, USA, 2010.

Jennifer Foster, Ozlem Cetinoglu, Joachim Wagner, Joseph Le Roux, Joakim Nivre, Deirdre Hogan, and Josef van Genabith. From news to comment: Resources and benchmarks for parsing the language of web 2.0. In *Proc. of the 5th International Joint Conference on Natural Language Processing (IJCNLP 2011)*, pages 893–901, Chiang Mai, Thailand, 2011. URL http://www.aclweb.org/anthology/I11-1100.

Kevin Gimpel, Nathan Schneider, Brendan O'Connor, Dipanjan Das, Daniel Mills, Jacob Eisenstein, Michael Heilman, Dani Yogatama, Jeffrey Flanigan, and Noah A. Smith. Part-of-speech tagging for Twitter: Annotation, features, and experiments. In *Proc. of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies (ACL HLT 2011)*, pages 42–47, Portland, USA, 2011. URL http://www.aclweb.org/anthology/P11-2008.

Stephan Gouws, Dirk Hovy, and Donald Metzler. Unsupervised mining of lexical variants from noisy text. In *Proceedings of the First workshop on Unsupervised Learning in NLP*, pages 82–90, Edinburgh, UK, 2011.

# References II

Bo Han and Timothy Baldwin. Lexical normalisation of short text messages: Makn sens a #twitter. In *Proc. of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies (ACL HLT 2011)*, pages 368–378, Portland, USA, 2011.

Bo Han, Paul Cook, and Timothy Baldwin. Automatically constructing a normalisation dictionary for microblogs. In *Proc. of the Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning 2012 (EMNLP-CoNLL 2012)*, pages 421–432, Jeju, Korea, 2012.

Fei Liu, Fuliang Weng, Bingqing Wang, and Yang Liu. Insertion, deletion, or substitution? normalizing text messages without pre-categorization nor supervision. In *Proc. of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies (ACL HLT 2011)*, pages 71–76, Portland, USA, 2011.

Fei Liu, Fuliang Weng, and Xiao Jiang. A broad-coverage normalization system for social media language. In *Proc. of the 50th Annual Meeting of the Association for Computational Linguistics (ACL 2012)*, pages 1035–1044, Jeju, Republic of Korea, 2012.

Marco Lui and Timothy Baldwin. Cross-domain feature selection for language identification. In *Proc. of the 5th International Joint Conference on Natural Language Processing (IJCNLP 2011)*, pages 553–561, Chiang Mai, Thailand, 2011.

Marco Lui and Timothy Baldwin. langid.py: An off-the-shelf language identification tool. In *Proc. of the 50th Annual Meeting of the Association for Computational Linguistics (ACL 2012) Demo Session*, pages 25–30, Jeju, Republic of Korea, 2012.

# References III

Alan Ritter, Sam Clark, Mausam, and Oren Etzioni. Named entity recognition in tweets: An experimental study. In *Proc. of the 2011 Conference on Empirical Methods in Natural Language Processing (EMNLP 2011)*, pages 1524–1534, Edinburgh, UK, 2011.

Joshua Ritterman, Miles Osborne, and Ewan Klein. Using prediction markets and Twitter to predict a swine flu pandemic. In *Proceedings of the 1st International Workshop on Mining Social Media*, November 2009.

Takeshi Sakaki, Makoto Okazaki, and Yutaka Matsuo. Earthquake shakes Twitter users: real-time event detection by social sensors. In *Proc. of the 19th International Conference on the World Wide Web (WWW 2010)*, pages 851–860, Raleigh, USA, 2010. ISBN 978-1-60558-799-8. doi: http://doi.acm.org/10.1145/1772690.1772777. URL http://doi.acm.org/10.1145/1772690.1772777.