

NLP for Social Media: Language Identification II and Text Normalization

Pawan Goyal

CSE, IITKGP

August 3-4, 2016

LI: Supervised Approaches

Input

- A document d
- A fixed set of classes $C = \{c_1, c_2, \dots, c_n\}$
- A training set of m hand-labeled documents $(d_1, c_1), \dots, (d_m, c_m)$

LI: Supervised Approaches

Input

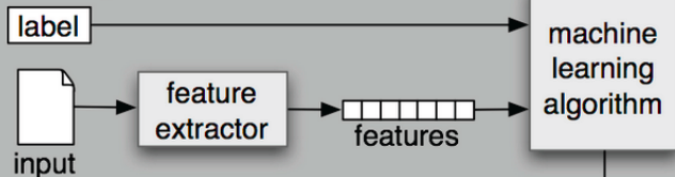
- A document d
- A fixed set of classes $C = \{c_1, c_2, \dots, c_n\}$
- A training set of m hand-labeled documents $(d_1, c_1), \dots, (d_m, c_m)$

Output

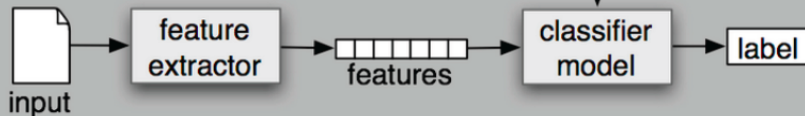
A learned classifier $\gamma: d \rightarrow c$

Supervised Machine Learning

(a) Training



(b) Prediction



Bayes' rule for documents and classes

For a document d and a class c

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

Bayes' rule for documents and classes

For a document d and a class c

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

Naïve Bayes Classifier

$$c_{MAP} = \arg \max_{c \in \mathcal{C}} P(c|d)$$

$$= \arg \max_{c \in \mathcal{C}} P(d|c)P(c)$$

$$= \arg \max_{c \in \mathcal{C}} P(x_1, x_2, \dots, x_n | c)P(c)$$

Naïve Bayes classification assumptions

$$P(x_1, x_2, \dots, x_n | c)$$

Naïve Bayes classification assumptions

$$P(x_1, x_2, \dots, x_n | c)$$

Bag of words assumption

Assume that the position of a word in the document doesn't matter

Naïve Bayes classification assumptions

$$P(x_1, x_2, \dots, x_n | c)$$

Bag of words assumption

Assume that the position of a word in the document doesn't matter

Conditional Independence

Assume the feature probabilities $P(x_i | c_j)$ are independent given the class c_j .

$$P(x_1, x_2, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \dots P(x_n | c)$$

Naïve Bayes classification assumptions

$$P(x_1, x_2, \dots, x_n | c)$$

Bag of words assumption

Assume that the position of a word in the document doesn't matter

Conditional Independence

Assume the feature probabilities $P(x_i | c_j)$ are independent given the class c_j .

$$P(x_1, x_2, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \dots P(x_n | c)$$

$$c_{NB} = \arg \max_{c \in C} P(c) \prod_{x \in X} P(x | c)$$

Maximum Likelihood Estimate

$$\hat{P}(c_j) = \frac{\text{doc-count}(C = c_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

Learning the model parameters

Maximum Likelihood Estimate

$$\hat{P}(c_j) = \frac{\text{doc-count}(C = c_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i|c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

Problem with MLE

Suppose in the training data, we haven't seen one of the words (say *pure*) in a given language.

$$\hat{P}(\text{pure}|\text{Hindi}) = 0$$

Learning the model parameters

Maximum Likelihood Estimate

$$\hat{P}(c_j) = \frac{\text{doc-count}(C = c_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

Problem with MLE

Suppose in the training data, we haven't seen one of the words (say *pure*) in a given language.

$$\hat{P}(\text{pure} | \text{Hindi}) = 0$$

$$c_{NB} = \arg \max_c \hat{P}(c) \prod_{x \in X} \hat{P}(x_i | c)$$

$$\begin{aligned}\hat{P}(w_i|c) &= \frac{\text{count}(w_i, c) + 1}{\sum_{w \in V} (\text{count}(w, c) + 1)} \\ &= \frac{\text{count}(w_i, c) + 1}{\left(\sum_{w \in V} \text{count}(w, c)\right) + |V|}\end{aligned}$$

A worked out example

	Doc	Words	Class
Training	1	English Wikipedia editor	en
	2	free English Wikipedia	en
	3	Wikipedia editor	en
	4	español de Wikipedia	es
Test	5	Wikipedia español el	?

A worked out example: No smoothing

	Doc	Words	Class
Training	1	English Wikipedia editor	en
	2	free English Wikipedia	en
	3	Wikipedia editor	en
	4	español de Wikipedia	es
Test	5	Wikipedia español el	?

$$\hat{P}(c) = \frac{\text{count}(c)}{\sum_{c_j \in \mathcal{C}} \text{count}(c_j)}$$

$$P(en) = 3/4 \quad P(sp) = 1/4$$

$$\hat{P}(t | c) = \frac{\text{count}(t, c)}{\sum_{t_i \in V} \text{count}(t_i, c)}$$

$$P(\text{"Wikipedia"} | en) = 3/8, \quad P(\text{"Wikipedia"} | es) = 1/3$$

$$P(\text{"español"} | en) = 0/8, \quad P(\text{"español"} | es) = 1/3$$

$$P(\text{"el"} | en) = 0/8, \quad P(\text{"el"} | es) = 0/3$$

$$P(en | doc5) = 3/4 \times 3/8 \times 0/8 \times 0/8 = 0$$

$$P(es | doc5) = 1/4 \times 2/9 \times 1/3 \times 0/3 = 0$$

A worked out example: Smoothing

	Doc	Words	Class
Training	1	English Wikipedia editor	en
	2	free English Wikipedia	en
	3	Wikipedia editor	en
	4	español de Wikipedia	sp
Test	5	Wikipedia español el	?

$$\hat{P}(c) = \frac{\text{count}(c)}{\sum_{c_j \in \mathcal{C}} \text{count}(c_j)}$$

$$P(en) = 3/4 \quad P(sp) = 1/4$$

$$\hat{P}(t | c) = \frac{\text{count}(t, c)}{\sum_{t_i \in V} \text{count}(t_i, c)}$$

$$P(\text{"Wikipedia"} | en) = 3 + 1/8 + 6, \quad P(\text{"Wikipedia"} | sp) = 1 + 1/3 + 6$$

$$P(\text{"español"} | en) = 0 + 1/8 + 6, \quad P(\text{"español"} | sp) = 1 + 1/3 + 6$$

$$P(\text{"el"} | en) = 0 + 1/8 + 6, \quad P(\text{"el"} | sp) = 0 + 1/3 + 6$$

$$P(en | doc5) = 3/4 \times 4/14 \times 1/14 \times 1/14 = 0.00109$$

$$P(sp | doc5) = 1/4 \times 2/9 \times 2/9 \times 1/9 = 0.00137$$

Character n -gram based Approach

Input: A word w (e.g., *kiprata*)

Character n -gram based Approach

Input: A word w (e.g., *kiprata*)

Features: character n -grams ($n=2$ to 5)

Character n -gram based Approach

Input: A word w (e.g., *khiprata*)

kshiprata \rightarrow \$kshiprata\$

2: \$k, ks, sh, hi, ip, pr, ra, at, ta, a\$

3: \$ks, ksh, shi, hip, ipr, ... ta\$

4: \$ksh, kshi, ship, ..., ata\$

5: \$kshi, kship, shipr, ..., rata\$

Features: character n -grams ($n=2$ to 5)

Character n -gram based Approach

Input: A word w (e.g., *khiprata*)

kshiprata \rightarrow \$kshiprata\$

2: \$k, ks, sh, hi, ip, pr, ra, at, ta, a\$

3: \$ks, ksh, shi, hip, ipr, ... ta\$

4: \$ksh, kshi, ship, ..., ata\$

5: \$kshi, kship, shipr, ..., rata\$

Features: character n -grams ($n=2$ to 5)

Classifier: Naïve Bayes, Max-Ent, SVMs

Character n -gram based Approach

Input: A word w (e.g., *khiprata*)

kshiprata \rightarrow \$kshiprata\$

2: \$k, ks, sh, hi, ip, pr, ra, at, ta, a\$

3: \$ks, ksh, shi, hip, ipr, ... ta\$

4: \$ksh, kshi, ship, ..., ata\$

5: \$kshi, kship, shipr, ..., rata\$

Features: character n -grams ($n=2$ to 5)

Classifier: Naïve Bayes, Max-Ent, SVMs

Prob (*kshiprata* is Sanskrit) \gg Prob (*kshiprata* is English)

langid.py	Lui and Baldwin [2012]
ChromeCLD	McCandless [2010]
LangDetect	Nakatani [2010]
LDIG	Nakatani [2012]
whatlang	Brown [2013]
YALI	Majliš [2012]
TextCat	Scheelen [2003]
MSR-LID	Goldszmidt et al. [2013]

```
python
Python 2.7.2+ (default, Oct 4 2011, 20:06:09)
[GCC 4.6.1] on linux2
Type "help", "copyright", "credits" or "license" for more information.
>>> import langid
>>> langid.classify("I do not speak english")
('en', 0.5713348767990674)
>>> langid.set_languages(['de','fr','it'])
>>> langid.classify("I do not speak english")
('it', 0.99999835791478453)
>>> langid.set_languages(['en','it'])
>>> langid.classify("I do not speak english")
('en', 0.99176190378750373)
```

<https://github.com/saffsd/langid.py>
Supports 97 languages

Word-level Language Labeling

Modi ke speech se India inspired ho gaya #namo

NE	Hn	En	Hn	NE	En	Hn	Hn	Other
	के		से			हो	गया	

Word-level Language Labeling

Modi ke speech se India inspired ho gaya #namo

NE	Hn	En	Hn	NE	En	Hn	Hn	Other
	के		से			हो	गया	

Modeling as a Sequence Prediction Problem

Given \mathbf{X} : $X_1 = \text{Modi}, X_2 = \text{ke}, \dots$

Word-level Language Labeling

Modi ke speech se India inspired ho gaya #namo

NE	Hn	En	Hn	NE	En	Hn	Hn	Other
	के		से			हो	गया	

Modeling as a Sequence Prediction Problem

Given \mathbf{X} : $X_1 = \text{Modi}, X_2 = \text{ke}, \dots$

Output: $\mathbf{Y} = Y_1$ (label for X_1), Y_2 (label for X_2), \dots

Word-level Language Labeling

Modi ke speech se India inspired ho gaya #namo

NE	Hn	En	Hn	NE	En	Hn	Hn	Other
	के		से			हो	गया	

Modeling as a Sequence Prediction Problem

Given \mathbf{X} : $X_1 = \text{Modi}, X_2 = \text{ke}, \dots$

Output: $\mathbf{Y} = Y_1$ (label for X_1), Y_2 (label for X_2), \dots

Such that $p(Y|X)$ is maximized

Conditional Random Fields: Modelling the Conditional Distribution

Conditional Random Fields: Modelling the Conditional Distribution

Model the Conditional Distribution:

$$P(\mathbf{y} | \mathbf{x})$$

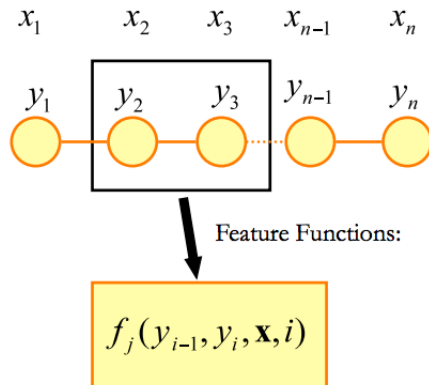
To predict a sequence compute:

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} P(\mathbf{y} | \mathbf{x})$$



Must be able to compute it efficiently.

Conditional Random Fields: Feature Functions



Express some characteristic of the empirical distribution that we wish to hold in the model distribution

$f_j(y_{i-1}, y_i, \mathbf{x}, i)$

1 if $y_{i-1} = IN$ and
 $y_i = NNP$ and
 $x_i = September$

0 otherwise

Label sequence modelled as a normalized product of feature functions:

$$P(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{Z(\mathbf{x})} \exp \sum_{i=1}^n \sum_j \lambda_j f_j(y_{i-1}, y_i, \mathbf{x}, i)$$

$$Z(\mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}} \sum_{i=1}^n \sum_j \lambda_j f_j(y_{i-1}, y_i, \mathbf{x}, i)$$

Features for word level Language Identification

Token-based features

- Capitalization
- Script
- Special Characters
- Character n-gram based classifiers
- Word length

Lexical Features

- Regular lexicon
- Unigram Frequency
- Entity Lexicon
- Acronym/slang lexicon

Context Features

- Next 3 tokens
- Last 3 tokens
- Current token
- Previous label (Bigram or B)

Characteristics of Text in Social Media

Social media text contains varying levels of “noise” (lexical, syntactic and otherwise), e.g.

Characteristics of Text in Social Media

Social media text contains varying levels of “noise” (lexical, syntactic and otherwise), e.g.

- Tell ppl u luv them cuz 2morrow is truly not promised.
- SUPER BOWL SUNDAY!!! Enjoy yourselves!!! Sunday morning GOODIES R sent out! C U 2Nyt!
- Follow @OFA today for more coverage of the gun violence petition delivery to Congress. #NotBackingDown #EarlyFF

Why is Social Media Text “Bad”?

Eisenstein [2013] identified the following possible contributing factors to “badness” in social media text:

- Lack of literacy?

Why is Social Media Text “Bad”?

Eisenstein [2013] identified the following possible contributing factors to “badness” in social media text:

- Lack of literacy? *no*

Why is Social Media Text “Bad”?

Eisenstein [2013] identified the following possible contributing factors to “badness” in social media text:

- Lack of literacy? *no*
- Length restrictions?

Why is Social Media Text “Bad”?

Eisenstein [2013] identified the following possible contributing factors to “badness” in social media text:

- Lack of literacy? *no*
- Length restrictions? *not primarily*

Why is Social Media Text “Bad”?

Eisenstein [2013] identified the following possible contributing factors to “badness” in social media text:

- Lack of literacy? *no*
- Length restrictions? *not primarily*
- Text input method-driven?

Why is Social Media Text “Bad”?

Eisenstein [2013] identified the following possible contributing factors to “badness” in social media text:

- Lack of literacy? *no*
- Length restrictions? *not primarily*
- Text input method-driven? *to some degree, yes*

Why is Social Media Text “Bad”?

Eisenstein [2013] identified the following possible contributing factors to “badness” in social media text:

- Lack of literacy? *no*
- Length restrictions? *not primarily*
- Text input method-driven? *to some degree, yes*
- Pragmatics (mimicking prosodic effects etc. in speech)?

Why is Social Media Text “Bad”?

Eisenstein [2013] identified the following possible contributing factors to “badness” in social media text:

- Lack of literacy? *no*
- Length restrictions? *not primarily*
- Text input method-driven? *to some degree, yes*
- Pragmatics (mimicking prosodic effects etc. in speech)? *yeeees*

Eisenstein, *What to do about bad language on the internet*, NAACL-HLT, 2013

What can be done about it?

Lexical normalization

Translate expressions into their canonical form

What can be done about it?

Lexical normalization

Translate expressions into their canonical form

Issues

- What are the candidate tokens for normalization?

What can be done about it?

Lexical normalization

Translate expressions into their canonical form

Issues

- What are the candidate tokens for normalization?
- To what degree do we allow normalization?

What can be done about it?

Lexical normalization

Translate expressions into their canonical form

Issues

- What are the candidate tokens for normalization?
- To what degree do we allow normalization?
- What is the canonical form of a given expression? (e.g., *aint*)

What can be done about it?

Lexical normalization

Translate expressions into their canonical form

Issues

- What are the candidate tokens for normalization?
- To what degree do we allow normalization?
- What is the canonical form of a given expression? (e.g., *aint*)
- Is lexical normalization always appropriate? (e.g., *bro*)

One standard definition

- relative to some standard tokenization

Task Definition

One standard definition

- relative to some standard tokenization
- consider only OOV tokens as candidates for normalization

Task Definition

One standard definition

- relative to some standard tokenization
- consider only OOV tokens as candidates for normalization
- allow only one-to-one word substitutions

i left ACL cus im sickk ! Yuu better be their tmrw
 ↓ ↓ ↓ ↓ ↓
i left ACL because I'm sick ! You better be their tomorrow

Task Definition

One standard definition

- relative to some standard tokenization
- consider only OOV tokens as candidates for normalization
- allow only one-to-one word substitutions

i left ACL cus im sickk ! Yuu better be their tmrw
 ↓ ↓ ↓ ↓ ↓
i left ACL because I'm sick ! You better be their tomorrow

Assumptions/corrolaries of the task definition:

- not possible to normalize in-vocabulary tokens, e.g. *their*

Task Definition

One standard definition

- relative to some standard tokenization
- consider only OOV tokens as candidates for normalization
- allow only one-to-one word substitutions

i left ACL cus im sickk ! Yuu better be their tmrw
 ↓ ↓ ↓ ↓ ↓
i left ACL because I'm sick ! You better be their tomorrow

Assumptions/corrolaries of the task definition:

- not possible to normalize in-vocabulary tokens, e.g. *their*
- not possible to normalise the multiword tokens, e.g., *ttyl*

Task Definition

One standard definition

- relative to some standard tokenization
- consider only OOV tokens as candidates for normalization
- allow only one-to-one word substitutions

i left ACL cus im sickk ! Yuu better be their tmrw
 ↓ ↓ ↓ ↓ ↓
i left ACL because I'm sick ! You better be their tomorrow

Assumptions/corrolaries of the task definition:

- not possible to normalize in-vocabulary tokens, e.g. *their*
- not possible to normalise the multiword tokens, e.g., *ttyl*
- ignore Twitter-specific entities, e.g., *obama*, *#mandela*, *bit.ly/1iRqm*

Task Definition

One standard definition

- relative to some standard tokenization
- consider only OOV tokens as candidates for normalization
- allow only one-to-one word substitutions

i left ACL cus im sickk ! Yuu better be their tmrw
 ↓ ↓ ↓ ↓ ↓
i left ACL because I'm sick ! You better be their tomorrow

Assumptions/corrolaries of the task definition:

- not possible to normalize in-vocabulary tokens, e.g. *their*
- not possible to normalise the multiword tokens, e.g., *ttyl*
- ignore Twitter-specific entities, e.g., *obama*, *#mandela*, *bit.ly/1iRqm*
- assume a unique correct “norm” for each token

Spelling Errors

TOMORROW

- Tomorow
- Tommorow
- Tommorrow

Phonetic/Cognitive
Errors

- Tpmorrow
- Tomrorow
- Tmorrow
- Tomnorrow

Typos or “slip of finger”
errors

Unintentional
Errors

Understanding unintentional spelling errors

TOMORROW

• Tomorow

Double letter omission

• Tommorow

Doubling of wrong letter

• Tommorrow

Doubling of letter

• Tpmorrow

Substitution: o \rightarrow p

• Tomrorow

Metathesis: or \rightarrow ro

• Tmorrow

Deletion: o \rightarrow ϵ

• Tomnorrow

Insertion: $\epsilon \rightarrow$ n

Phonetic/Cognitive
Errors

Typos or "slip of finger"
errors

- Cost of **Edit Operations**:

- Insertion($\epsilon \rightarrow c$): 1
- Deletion ($c \rightarrow \epsilon$): 1
- Substitution: ($c \rightarrow c'$): 1 or 2

Metathesis ($cc' \rightarrow c'c$) is either modeled as a single edit operation (cost = 1) or as a deletion-insertion pair ($cc' \rightarrow \epsilon c' \rightarrow c'c$), and therefore cost of 2.

- **Edit Distance** between two strings $s:c_1c_2c_3\dots c_n$ and $s':c'_1c'_2c'_3\dots c'_n$ is defined as the minimum value of the sum of the cost of a sequence of edit operations required to convert s to s' .

- *engine & begin*, *elevator & evaluator*, *east & csar*

- Dynamic Programming Algorithm

What about spelling errors in Social Media?

The shorter → the faster
Constraint: understandability

24

dis is n eg 4 txtin lang

39

This is an example for Texting language

Other factors: Coolness, group-membership, accommodating

The case of 'Tomorrow'

- 2moro (9)
- tomoz (25)
- tomoro (12)
- tomrw (5)
- tom (2)
- tomra (2)
- tomorrow (24)
- tomora (4)
- tomm (1)
- tomo (3)
- tomorrow (3)
- 2mro (2)
- morrow (1)
- tomor (2)
- tmorro (1)
- moro (1)

Spell-checkers, such as Aspell, perform very poorly on such data (<22%)

Data from (Choudhury et al., 2007)

Patterns or Compression Operators

Phonetic substitution (phoneme)

Patterns or Compression Operators

Phonetic substitution (phoneme)

psycho → syco, then → den

Patterns or Compression Operators

Phonetic substitution (phoneme)

psycho → syco, then → den

Phonetic substitution (syllable)

Patterns or Compression Operators

Phonetic substitution (phoneme)

psycho → syco, then → den

Phonetic substitution (syllable)

today → 2day, see → c

Patterns or Compression Operators

Phonetic substitution (phoneme)

psycho → syco, then → den

Phonetic substitution (syllable)

today → 2day, see → c

Deletion of vowels

Patterns or Compression Operators

Phonetic substitution (phoneme)

psycho → syco, then → den

Phonetic substitution (syllable)

today → 2day, see → c

Deletion of vowels

message → mssg, about → abt

Patterns or Compression Operators

Phonetic substitution (phoneme)

psycho → syco, then → den

Phonetic substitution (syllable)

today → 2day, see → c

Deletion of vowels

message → mssg, about → abt

Deletion of repeated characters

Patterns or Compression Operators

Phonetic substitution (phoneme)

psycho → syco, then → den

Phonetic substitution (syllable)

today → 2day, see → c

Deletion of vowels

message → mssg, about → abt

Deletion of repeated characters

tomorrow → tomorow

Truncation (deletion of tails)

Patterns or Compression Operators

Truncation (deletion of tails)

introduction → intro, evaluation → eval

Patterns or Compression Operators

Truncation (deletion of tails)

introduction → intro, evaluation → eval

Common Abbreviations

Patterns or Compression Operators

Truncation (deletion of tails)

introduction → intro, evaluation → eval

Common Abbreviations

Kharagpur → kgp, text back → tb

Patterns or Compression Operators

Truncation (deletion of tails)

introduction → intro, evaluation → eval

Common Abbreviations

Kharagpur → kgp, text back → tb

Informal pronunciation

Patterns or Compression Operators

Truncation (deletion of tails)

introduction → intro, evaluation → eval

Common Abbreviations

Kharagpur → kgp, text back → tb

Informal pronunciation

going to → gonna

Patterns or Compression Operators

Truncation (deletion of tails)

introduction → intro, evaluation → eval

Common Abbreviations

Kharagpur → kgp, text back → tb

Informal pronunciation

going to → gonna

Emphasis by repetition

Patterns or Compression Operators

Truncation (deletion of tails)

introduction → intro, evaluation → eval

Common Abbreviations

Kharagpur → kgp, text back → tb

Informal pronunciation

going to → gonna

Emphasis by repetition

Funny → fuuuunnnnnnyyyyyy

Successive Application of Operators

- Because \rightarrow cause (informal usage)

Successive Application of Operators

- Because → cause (informal usage)
- cause → cauz (phonetic substitution)

Successive Application of Operators

- Because → cause (informal usage)
- cause → cauz (phonetic substitution)
- cauz → cuz (vowel deletion)

Categorisation of non-standard words in English Twitter

Category	Ratio	Example
Letter&Number	2.36%	<i>b4</i> "before"
Letter	72.44%	<i>shuld</i> "should"
Number Substitution	2.76%	<i>4</i> "for"
Slang	12.20%	<i>lol</i> "laugh out loud"
Other	10.24%	<i>sucha</i> "such a"

Table : Types of non-standard words in a 449 message sample of English tweets

Categorisation of non-standard words in English Twitter

Category	Ratio	Example
Letter&Number	2.36%	<i>b4</i> "before"
Letter	72.44%	<i>shuld</i> "should"
Number Substitution	2.76%	<i>4</i> "for"
Slang	12.20%	<i>lol</i> "laugh out loud"
Other	10.24%	<i>sucha</i> "such a"

Table : Types of non-standard words in a 449 message sample of English tweets

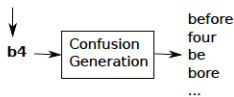
Most non-standard words in sampled tweets are *morphophonemic* variations.

Token-based Approach (Han and Baldwin, 2011)

Token-based Approach (Han and Baldwin, 2011)

- 1 Confusion set generation (i.e., find correction candidates)
- 2 Non-standard word detection (i.e., is the OOV a non-standard word?)
- 3 Normalisation of a non-standard word (i.e., select the candidate)

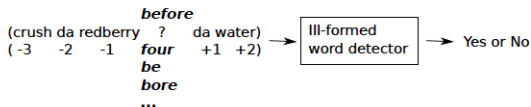
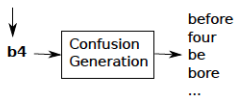
... crush da redberry **b4** da water ...



Token-based Approach (Han and Baldwin, 2011)

- 1 Confusion set generation (i.e., find correction candidates)
- 2 Non-standard word detection (i.e., is the OOV a non-standard word?)
- 3 Normalisation of a non-standard word (i.e., select the candidate)

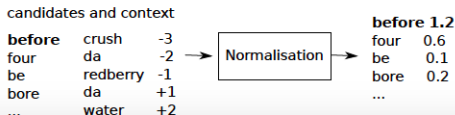
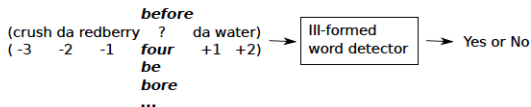
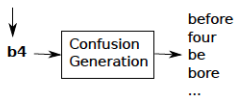
... crush da redberry **b4** da water ...



Token-based Approach (Han and Baldwin, 2011)

- 1 Confusion set generation (i.e., find correction candidates)
- 2 Non-standard word detection (i.e., is the OOV a non-standard word?)
- 3 Normalisation of a non-standard word (i.e., select the candidate)

... crush da redberry **b4** da water ...



- Filter out any Twitter-specific tokens (user-mentions, hashtags, RT, etc.) and URLs

- Filter out any Twitter-specific tokens (user-mentions, hashtags, RT, etc.) and URLs
- Identify all OOV words relative to a standard spelling dictionary (aspell)

- Filter out any Twitter-specific tokens (user-mentions, hashtags, RT, etc.) and URLs
- Identify all OOV words relative to a standard spelling dictionary (aspell)
- For OOV words, shorten any repetitions of 3+ letters to 2 letters

- Generation via edit distance over letters (T_c) and phonemes (T_p).

Candidate Generation

- Generation via edit distance over letters (T_c) and phonemes (T_p).
- This allows to generate “earthquake” for words such as *earthquick*.

- Generation via edit distance over letters (T_c) and phonemes (T_p).
- This allows to generate “earthquake” for words such as *earthquick*.
- Candidates with $T_c \leq 2 \vee T_p \leq 1$ were taken, further filtered using frequency to take the top 10% of candidates.

Detection of Ill-formed words

Detection based on candidate context fitness

- Correct words should fit better with context than substitution candidates
- Incorrect words should fit *worse* than substitution candidates

Detection of Ill-formed words

Detection based on candidate context fitness

- Correct words should fit better with context than substitution candidates
- Incorrect words should fit *worse* than substitution candidates

Basic Idea: Use Dependencies from corpus data

An SVM classifier is trained based on dependencies, to indicate candidate context fitness.

Ill-formed word in text snippet	Candidate	Dependencies
<i>but I was thinkin movies .</i>	(thinking, ...)	dobj(thinking, movies)
<i>article poster by runderobb : there was</i>	(ratrap, ...)	-

Feature Representation using Dependencies

- Build a dependency bank from existing corpora
- Represent each dependency tuple as a word pair + positional index

Feature Representation using Dependencies

- Build a dependency bank from existing corpora
- Represent each dependency tuple as a word pair + positional index

Corpus (NYT)

One obvious difference is the way they look, ...



Stanford Parser

num(difference-3, One-1)
amod(difference-3, obvious-2)
nsubj(way-6, difference-3)
cop(way-6, is-4)
det(way-6, the-5)
dobj(look-8, way-6)
nsubj(look-8, they-7)
rcmod(way-6, look-8)
...



Dependency bank

(way, difference, 3)
(look, way, 2)
...

SVM Training Data Generation

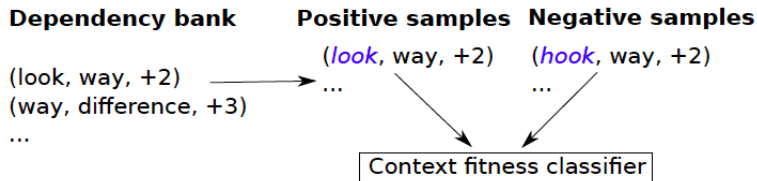
- Use dependency bank directly as positive features

SVM Training Data Generation

- Use dependency bank directly as positive features
- Automatically generate negative dependency features by replacing the target word with highly-ranked confusion candidates

SVM Training Data Generation

- Use dependency bank directly as positive features
- Automatically generate negative dependency features by replacing the target word with highly-ranked confusion candidates



Detecting ill-formed words

OOV words with candidates fitting the context (i.e., positive classification outputs) are probably ill-formed words

Detecting ill-formed words

OOV words with candidates fitting the context (i.e., positive classification outputs) are probably ill-formed words

...way yu **lookin** shuld be a sin ..

↓
looking ~ (way, looking, -2)
hooking ~ (way, hooking, -2)

...

*If positive outputs exceed the threshold,
feed all candidates for normaliation.*

Prediction
Context fitness classifier

↓
looking ~ +1
hooking ~ -1

...

Detecting ill-formed words

OOV words with candidates fitting the context (i.e., positive classification outputs) are probably ill-formed words

...way yu **lookin** shuld be a sin ..

↓
looking ~ (way, looking, -2)

hooking ~ (way, hooking, -2)

...

*If positive outputs exceed the threshold,
feed all candidates for normaliation.*

Prediction
Context fitness classifier

↓
looking ~ +1
hooking ~ -1

...

Threshold = 1 → *lookin* is considered to be an ill-formed word

For each ill-formed word and its possible correction candidates, the following features are considered for normalization:

Word Similarity

- letter and phoneme edit distance (ED)
- prefix, suffix, and longest common subsequence

Normalization Candidate Selection

For each ill-formed word and its possible correction candidates, the following features are considered for normalization:

Word Similarity

- letter and phoneme edit distance (ED)
- prefix, suffix, and longest common subsequence

Context Support

- trigram language model score
- dependency score (weighted dependency count, derived from the detection step)

Type-based approach

Observation

The longer the ill-formed word, the more likely there is a unique normalization candidate

Observation

The longer the ill-formed word, the more likely there is a unique normalization candidate

- $y \Rightarrow \{\underline{why}, \underline{you}, \dots\}$, $hw \Rightarrow \{\underline{how}, \underline{homework}, \dots\}$
- $4eva \Rightarrow \{\underline{forever}\}$, $tlkin \Rightarrow \{\underline{talking}\}$

Type-based approach

Observation

The longer the ill-formed word, the more likely there is a unique normalization candidate

- $\underline{y} \Rightarrow \{\underline{why}, \underline{you}, \dots\}$, $\underline{hw} \Rightarrow \{\underline{how}, \underline{homework}, \dots\}$
- $\underline{4eva} \Rightarrow \{\underline{forever}\}$, $\underline{tlnkin} \Rightarrow \{\underline{talking}\}$

Approach

Construct a dictionary of (lexical variant, standard form) pair for longer word types (character length ≥ 4) of moderate frequency (≥ 16)

Type-based Approach (Han et al. (2012))

Construct the dictionary based on distributional similarity + string similarity

Type-based Approach (Han et al. (2012))

Construct the dictionary based on distributional similarity + string similarity

Input: Tokenised English tweets

- Extract (OOV, IV) pairs based on distributional similarity
- Re-rank the extracted pairs by string similarity

Type-based Approach (Han et al. (2012))

Construct the dictionary based on distributional similarity + string similarity

Input: Tokenised English tweets

- Extract (OOV, IV) pairs based on distributional similarity
- Re-rank the extracted pairs by string similarity

Output

A list of (OOV, IV) pairs ordered by string similarity; select the top- n pairs for inclusion in the normalisation lexicon.

An Example

... see you *tmrw* ...
... *tmrw* morning ...
... *tomorrow* morning ...
...

⇓ distributional similarity

{*tmrw*, *2morow*, *tomorrow*, *Monday*}

⇓ string similarity

tmrw → *tomorrow*

Components/parameters of the method

- context window size: ± 1 , ± 2 , ± 3
- context word sensitivity: bag-of-words vs. positional indexing
- context word representation: unigram, bigram or trigram
- context word filtering: all tokens vs. only dictionary words
- context similarity: KL divergence, Jensen-Shannon divergence, Cosine similarity, Euclidean distance

Components/parameters of the method

- context window size: ± 1 , ± 2 , ± 3
- context word sensitivity: bag-of-words vs. positional indexing
- context word representation: unigram, bigram or trigram
- context word filtering: all tokens vs. only dictionary words
- context similarity: KL divergence, Jensen-Shannon divergence, Cosine similarity, Euclidean distance

Tune parameters relative to (OOV,IV) pair development data

Rerank pairs by string similarity

(OOV,IV) pairs derived by distributional similarity:

(*Obama, Adam*) ↓

(*tmrw, tomorrow*) ↑

(*Youtube, web*) ↓

(*4eva, forever*) ↑

...

Rerank pairs by string similarity

(OOV,IV) pairs derived by distributional similarity:

(*Obama, Adam*) ↓

(*tmrw, tomorrow*) ↑

(*Youtube, web*) ↓

(*4eva, forever*) ↑

...

(*tmrw, tomorrow*)

(*4eva, forever*)

Get the top-ranked pairs as lexicon entries:

- Han, Bo, and Timothy Baldwin. “Lexical normalisation of short text messages: Makn sens a# twitter.” Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1. Association for Computational Linguistics, 2011.
- Han, Bo, Paul Cook, and Timothy Baldwin. “Automatically constructing a normalisation dictionary for microblogs.” Proceedings of the 2012 joint conference on empirical methods in natural language processing and computational natural language learning. Association for Computational Linguistics, 2012.