Racial Disparity in Natural Language Processing: A Case Study of Social Media African-American English

Su Lin Blodgett



### Brendan O'Connor



Talk at FAT/ML Workshop -- August 2017 College of Information and Computer Sciences University of Massachusetts Amherst Paper, etc.: <u>http://brenocon.com</u>

# Disparity in language technology

- Language technologies analyze the linguistic behavior of people
- Language is affected by social context and attributes

1.00 0.75 Example: gender bias in Word Error Rate YouTube autocaptions [Tatman 2017, EACL Ethics in NLP workshop] 0.25

- 0.00 Women Men
- What information can a user access?
- Whose voices are heard?

- Implications of socially embedded language for natural language processing?
- Twitter POS tagging: Race matters!
- African-American English: Language ID and parsing
- Goal: have NLP tools work well across dialects and genres/mediums
  - Language technology to serve dialect speakers

## Kids these days



[Eisenstein, O'Connor, Smith, Xing, PLOS ONE 2014]

### NLP on social media's own terms

ikr	smh	he	asked	fir	уо	last
name	SO	he	can	add	u	on
fb	lololol					

- Is this "noisy text"?
- Any NLP system, starting with POS tagging, needs different models/resources than traditional written English
- Word clusters on unlabeled tweets (56 million)

[Owoputi et al. NAACL 2013, Gimpel et al. ACL 2011, www.cs.cmu.edu/~ark/TweetNLP]

Monday, August 14, 17

## NLP on social media's own terms



w fo fa fr fro ov fer fir whit abou aft serie fore fah fuh w/her w/that fron isn agains"non-standard<br/>prepositions"yeah yea nah naw yeahh nooo yeh noo noooo yeaa ikr nvm yeahhh nahh nooooo"interjections"facebook fb itunes myspace skype ebay tumblr bbm flickr aim msn netflix pandora"online service<br/>names"smh jk #fail #random #fact smfh #smh #winning #realtalk smdh #dead #justsaying"hashtag-y<br/>interjections"?

[Owoputi et al. NAACL 2013, Gimpel et al. ACL 2011, www.cs.cmu.edu/~ark/TweetNLP]

Monday, August 14, 17

# What does it learn?

Orthographic normalizations

so s0 -so so- \$o /so //so

# (Immediate?) future auxiliaries

gonna gunna gona gna guna gnna ganna qonna gonna gana qunna gonne goona gonnaa g0nna goina gonnah goingto gunnah gonaa gonan gunnna going2 gonnna gunnaa gonny gunaa quna goonna qona gonns goinna gonnae qnna gonnaaa gnaa

tryna gon finna bouta trynna boutta gne fina gonn tryina fenna qone trynaa qon boutaa funna finnah bouda boutah abouta fena bouttah boudda trinna qne finnaa fitna aboutta goin2 bout2 finnna trynah finaa ginna bouttaa fna try'na g0n trynn tyrna trna bouto finsta fnna tranna finta tryinna finnuh tryingto boutto

- finna ~ "fixing to"
- tryna ~ "trying to"
- bouta ~ "about to"

# Subject-AuxVerb constructs



• Where do nonstandard terms come from?



Q,

https://twitter.com/search?q=imma&src=typd&vertical=default&f=tweets





https://twitter.com/search?q=imma&src=typd&vertical=default&f=tweets



Demographic Dialectal Variation in Social Media: A Case Study of African-American English





#### Brendan O'Connor



### **EMNLP 2016**

### Dialect

### he woke af smart af educated af daddy af coconut oil using af GOALS AF & shares food af







### Why is social media different?

- Internet speech?
- Pre-existing dialectal English?
  - Geographic patterns of word usage often reveal relationships to race, ethnicity etc.
  - African-American English in Twitter [Eisenstein 2013, Jorgensen et al. 2015, Jones 2015]



### Youth, minorities on Twitter





### Youth, minorities on Twitter

#### [Pew Research]



P(use twitter | age)



Monday, August 14, 17

- Learn language models correlated with U.S. Census racial demographics
- Validate against sociolinguistic knowledge
- **Investigate** racial disparities in natural language processing tools
- Adapt/create fair NLP tools

### Associating geolocated tweets with demographics



2+ Follow

he woke af smart af educated af daddy af coconut oil using af GOALS AF & shares food af

Bored af den my phone finna die!!!

12

### Associating geolocated tweets with demographics

2+ Follow



he woke af smart af educated af daddy af

coconut oil using af GOALS AF & shares food af

Bored af den my phone finna die!!!



15



19

Monday, August 14, 17

### Corpus creation and linguistic validation

- Beyond unigrams: creation of user-level topic-aligned corpora
- How do we linguistically validate them?
  - Lexicon
  - Phonology (Jones, Jorgensen et al.)
  - Syntax (Stewart)

### Lexical analysis

• For every word in vocabulary w and topic k, calculate

$$r_k(w) = \frac{p(w|z=k)}{p(w|z\neq k)}$$

- Examine w where  $r_{AA}(w) \ge 2$ ,  $r_{white}(w) \ge 2$ : AA- and whitealigned words
- 79% of AA-aligned words, 58% of white-aligned words not in a standard English dictionary

### Phonological analysis

- Calculate r<sub>AA</sub>(w) for 31 phonological variants illustrated through nonstandard spellings
- For 30/31 variants: *r* ≥ 1

AAE	Ratio	SAE
sholl	1802.49	sure
iont	930.98	I don't
wea	870.45	where
talmbout	809.79	talking about
sumn	520.96	something

29

### Examining NLP tools - language identification

• Language identification - key step in NLP pipelines

```
>>> s = '''he woke af smart af educated af daddy af coconut oil using af
... GOALS AF & shares food af'''
>>> langid.classify(s)
('da', 0.9999999993212958)
>>> s = 'Bored af den my phone finna die!!!'
>>> langid.classify(s)
('da', 0.9999968001354156)
```

AA Acc

WH Acc

- p(correct | Wh) vs
   p(correct | AA)
- Assess disparity in
  - langid.py: popular
     open-source system
     [Lui and Baldwin, 2012]
  - Twitter (in metadata)
  - IBM, Microsoft

AA Acc

WH Acc

- p(correct | Wh) vs
   p(correct | AA)
- Assess disparity in
  - langid.py: popular
     open-source system
     [Lui and Baldwin, 2012]
  - Twitter (in metadata)
  - IBM, Microsoft



AA Acc

WH Acc

- p(correct | Wh) vs
   p(correct | AA)
- Assess disparity in
  - langid.py: popular
     open-source system
     [Lui and Baldwin, 2012]
  - Twitter (in metadata)
  - IBM, Microsoft



		AA Acc.	WH Acc.	Diff.
	$t \leq 5$	68.0	70.8	2.8
langid by	$5 < t \le 10$	84.6	91.6	7.0
iangia.py	$10 < t \le 15$	93.0	98.0	5.0
	<i>t</i> > 15	96.2	99.8	3.6
	$t \leq 5$	62.8	77.9	15.1
IBM Wataan	$5 < t \le 10$	91.9	95.7	3.8
ibivi watsoli	$10 < t \le 15$	96.4	99.0	2.6
	<i>t</i> > 15	98.0	99.6	1.6
	$t \leq 5$	87.6	94.2	6.6
Microsoft Agura	$5 < t \le 10$	98.5	99.6	1.1
MICIOSOIT AZUIE	$10 < t \le 15$	99.6	99.9	0.3
	<i>t</i> > 15	99.5	99.9	0.4
	$t \leq 5$	54.0	73.7	19.7
Twitton	$5 < t \le 10$	87.5	91.5	4.0
IWILLEI	$10 < t \le 15$	95.7	96.0	0.3
	<i>t</i> > 15	98.5	95.1	-3.0

# Less disparate language ID

- Our approach: use our model's demographic inferences as additional signal of English-ness: Allows for racial and online dialect
  - Effectively expands identifier's training set
- Improves English recall, even for international tweets (!) [Blodgett and O'Connor, WNUT 2017]
- Jurgens et al., ACL 2017: Variety of methods for broader training data also improve English recall

### Examining NLP tools - dependency parsing

- Compare annotated parses to systems' output parses
  - **WIRED** Google Has Open Sourced SyntaxNet, Its AI for Understanding Language



#### Announcing SyntaxNet: The World's Most Accurate Parser Goes Open Source

Thursday, May 12, 2016

Posted by Slav Petrov, Senior Staff Research Scientist



32

### Examining NLP tools - dependency parsing

- Compare annotated parses to systems' output parses
- AAE-like tweets are much harder than SAE-like tweets

Parser	AA	Wh.	Difference
SyntaxNet	64.0 (2.5)	80.4 (2.2)	16.3 (3.4)
CoreNLP	50.0 (2.7)	71.0 (2.5)	21.0 (3.7)

Recall for annotated edges for each message set, bootstrapped standard errors in parentheses.

34

Ongoing work: create a fair parser

Online language captures historical AAE variation (Uniformly sampled Twitter = lots of AAE!)

#### 1914: reported speech

(Elizabeth Waties Allston Pringle, "A Woman Rice Planter," First-Person Narratives of the American South Collection)

dey b'longs to dat gent'man ahaid

Online language captures historical AAE variation (Uniformly sampled Twitter = lots of AAE!)

#### 1914: reported speech

(Elizabeth Waties Allston Pringle, "A Woman Rice Planter," First-Person Narratives of the American South Collection)

#### dey b'longs to dat gent'man ahaid

2013: Twitter data

$$\frac{P(\det | AA)}{P(\det | \neg AA)} = 5.9$$

$$\frac{P(\text{dey} \mid AA)}{P(\text{dey} \mid \neg AA)} = 6.8$$

Online language captures historical AAE variation (Uniformly sampled Twitter = lots of AAE!)

#### 1914: reported speech

(Elizabeth Waties Allston Pringle, "A Woman Rice Planter," First-Person Narratives of the American South Collection)

#### dey b'longs to dat gent'man ahaid

#### 2013: Twitter data

$$\frac{P(\det | AA)}{P(\det | \neg AA)} = 5.9 \qquad \qquad \frac{P(\det | AA)}{P(\det | \neg AA)} = 6.8$$

#### POS taggers: standard vs. designed for Twitter

CoreNLP	<b>dey/NN(PRP)</b> b/NN(VBZ) '/Punct longs/NNS(VBZ) to/TO
	<b>dat</b> / <b>VB</b> ( <b>DT</b> ) gent/JJ '/Punct man/NN ahaid/ <b>VBN</b> ( <b>RB</b> )
ARK	<b>dey</b> /Pro b'longs/Verb to/Prep <b>dat</b> /Det gent'man/Noun ahaid/Adv

# NLP and social bias

- Natural language processing (NLP) resources are typically designed for standard English or other major languages
  - But non-standard languages correlates with social background
- How do social confounds affect other language technologies?
  - Sentiment measurement? Political science? Digital humanities?
  - Search? Translation?
- How to adapt NLP systems
- Online data from social processes reproduces social phenomena, and algorithms re-learn it