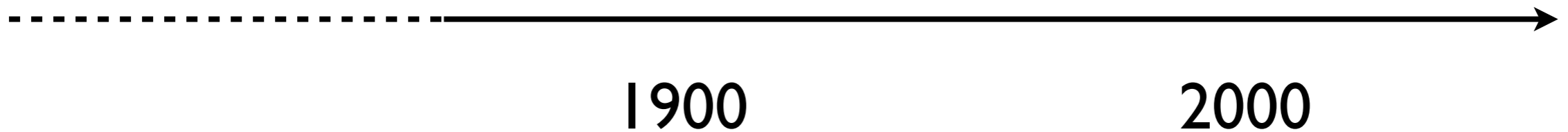


Demographic bias in social media  
language analysis:  
a case study of African-American English

Talk at Network Science Institute,  
Northeastern University, Dec. 14 2016

Brendan O'Connor (<http://brenocon.com>)  
College of Information and Computer Sciences  
University of Massachusetts Amherst

# Computational Social Science

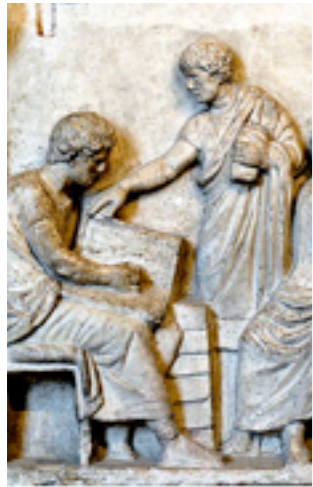


# Computational Social Science

## Official social data

---

Data collection    Data analysis    Data computation



100 BCE



1829



1890



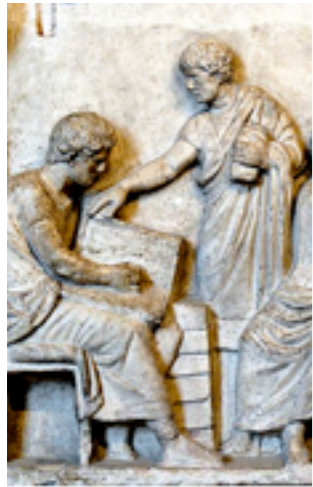
1900

2000

# Computational Social Science

## Official social data

Data collection    Data analysis    Data computation



100 BCE



1829



1890

## Semi-structured social data



### Digitized behavior

Billions of users,  
messages/day



### Digitized news

Thousands of articles/day



### Digitized archives

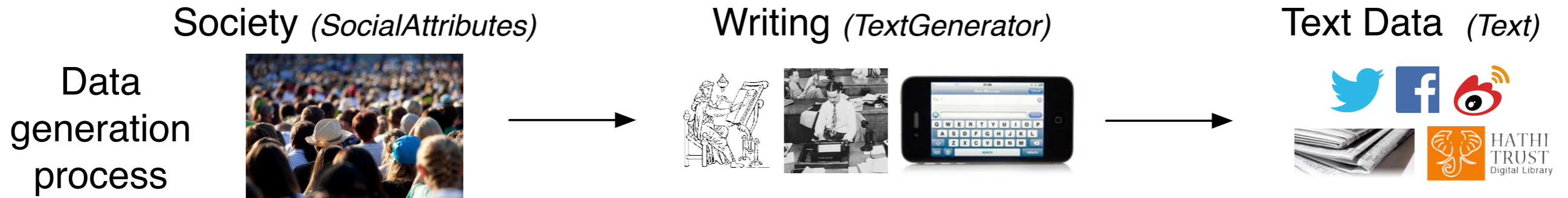
Millions of books/century



1900

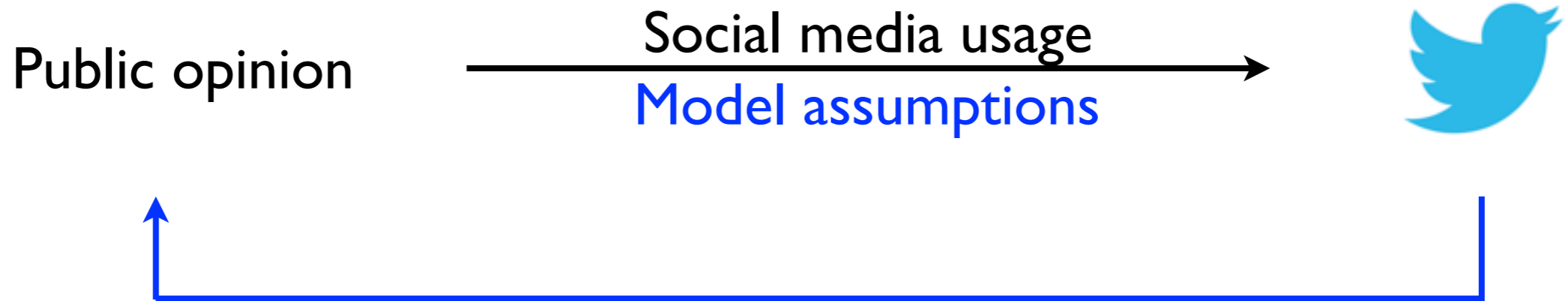
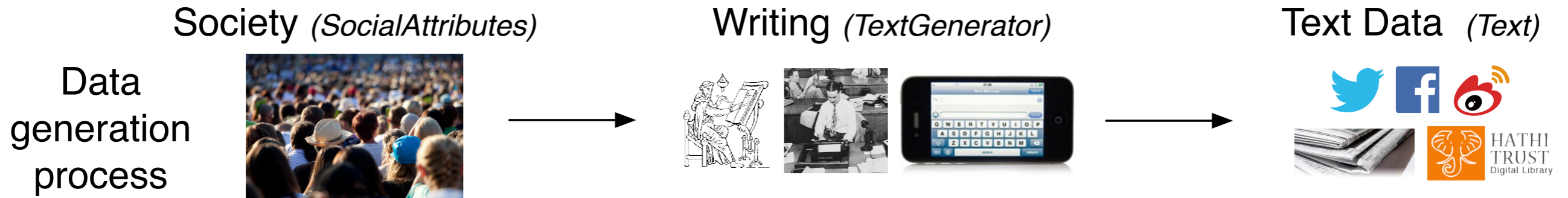
2000

# $TextGenerator(SocialAttributes) \rightarrow Text$

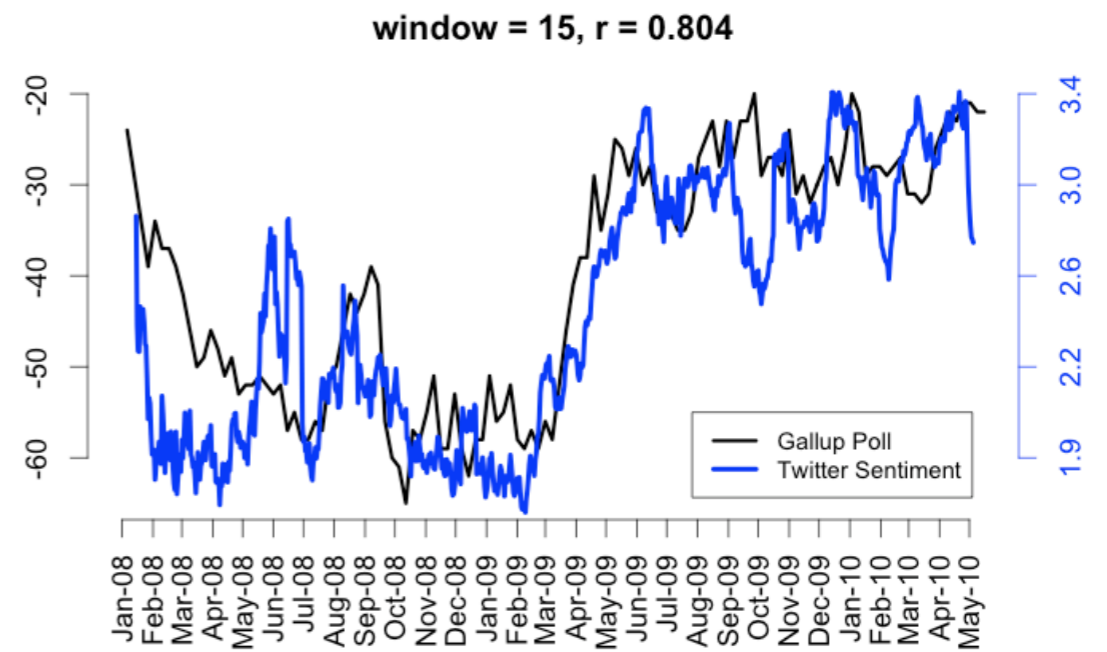


Language for social measurement  
 $P(\text{SocAttr} \mid \text{Text}, \text{TextGen})$

# TextGenerator(SocialAttributes) → Text

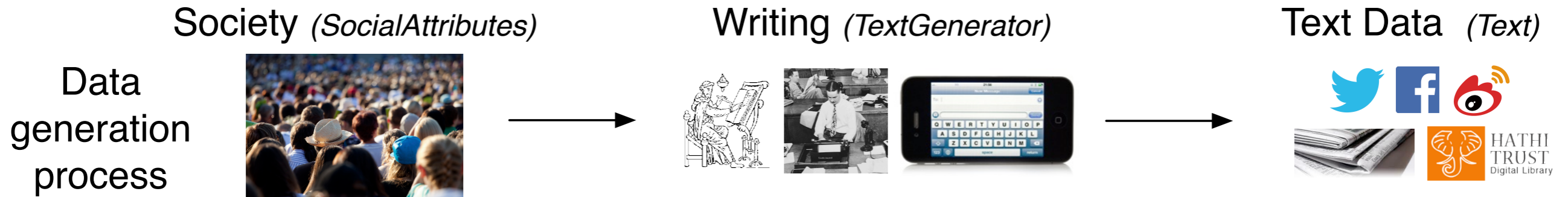


Language for social measurement  
 $P(\text{SocAttr} \mid \text{Text}, \text{TextGen})$



[O'Connor et al., ICWSM 2010]

# TextGenerator(SocialAttributes) → Text



Real-world political events

News media process  
Model assumptions

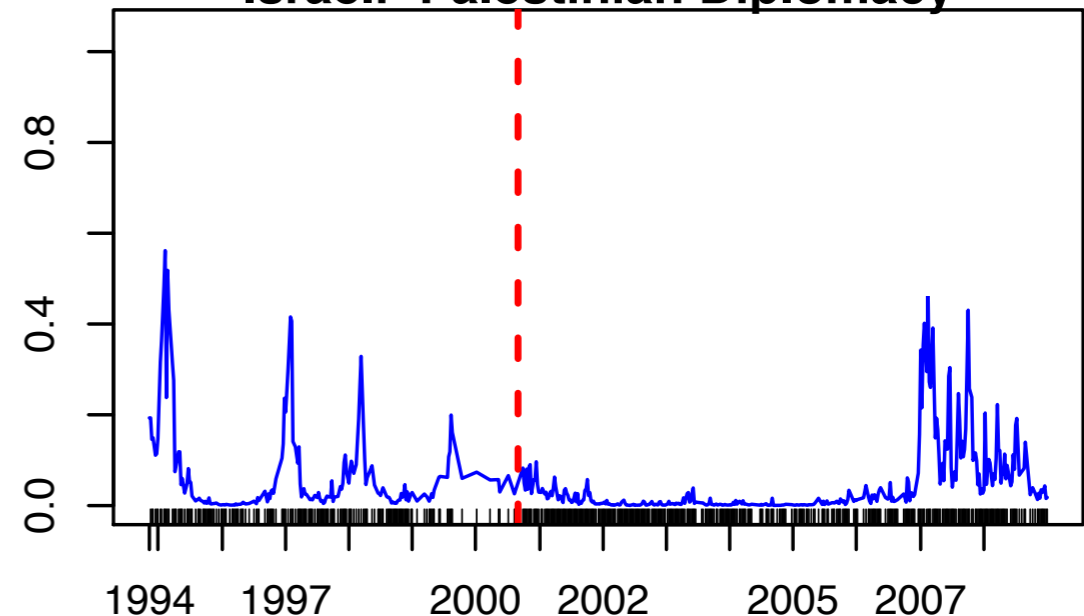


Language for social measurement  
 $P(\text{SocAttr} \mid \text{Text}, \text{TextGen})$

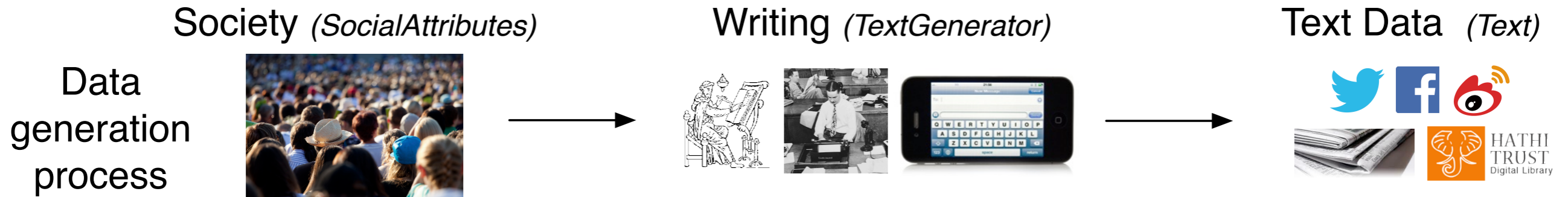
[meet with, sign with, praise, say with, arrive in, host, tell, welcome, join, thank, meet, travel to, criticize, leave, take to, begin to, begin with, summon, reach with, hold with...]

[O'Connor, Stewart, Smith ACL 2013]

Israeli-Palestinian Diplomacy



# *TextGenerator(SocialAttributes) → Text*



## What to analyze:

Social phenomena in social media datasets

- Political speech under Chinese censorship
- Events in international relations
- Social factors in language use

## How to analyze:

NLP capabilities we need to do these better

- Part of speech tagging
- Entity extraction
- Syntactic, semantic parsing

**What social bias exists in NLP models?**



# Linguistic/speech act diversity on Twitter

Official announcements



**BritishMonarchy** TheBritishMonarchy  
On 6 Jan: Changing the Guard at Buckingham Palace - Starts at approx 11am <http://www.royal.gov.uk/G>  
17 hours ago

Business advertising



**bigdogcoffee** bigdogcoffee  
Back to normal hours beginning tomorrow.....Monday-Friday 6am-10pm Sat/Sun 7:30am-10pm  
2 Jan

Links to blog and web content



**crampell** Catherine Rampell  
Casey B. Mulligan: Assessing the Housing Sector - <http://nyti.ms/hcUKK9>  
10 hours ago

Celebrity self-promotion



**THE\_REAL\_SHAQ** THE\_REAL\_SHAQ  
fill in da blank, my new years shaqalution is \_\_\_\_\_  
4 Jan

Status messages



**emax** electronic max  
1.1.11 - britons and americans can agree on the date for once. happy binary day!  
1 Jan

Group conversation



**\_siddx3** Evelyn Santana  
RT @\_LusciousVee: [#EveryoneShouldKnow](#) Ima Finally Be 18 This Year ^^  
3 minutes ago

Personal conversation



**xoxoJuicyCee** CeeCee♥  
[@fxknnCelly](#) aha kayy goodnightt (:  
4 Jan

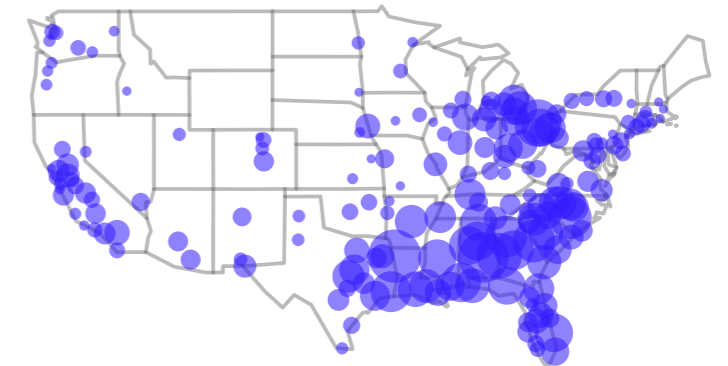
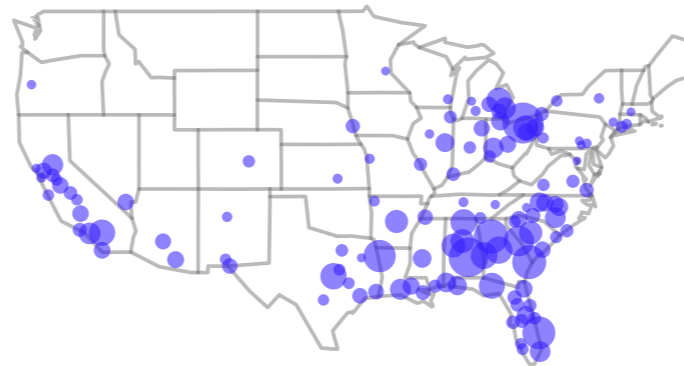
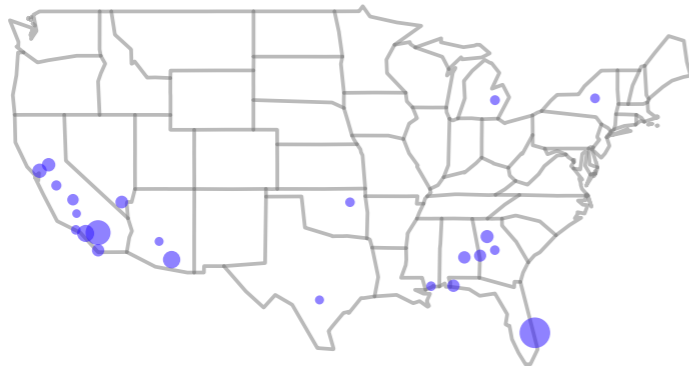
# Kids these days

weeks 1–50

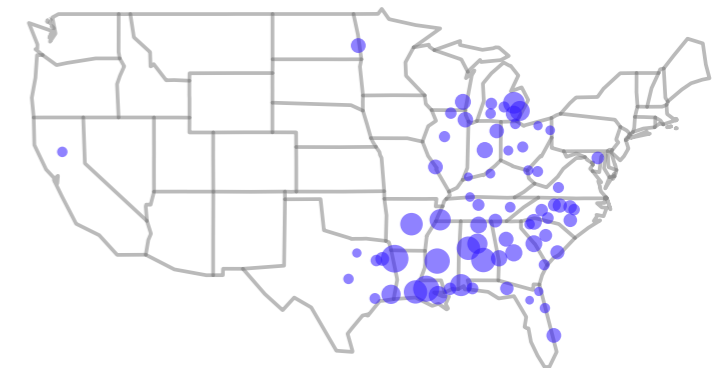
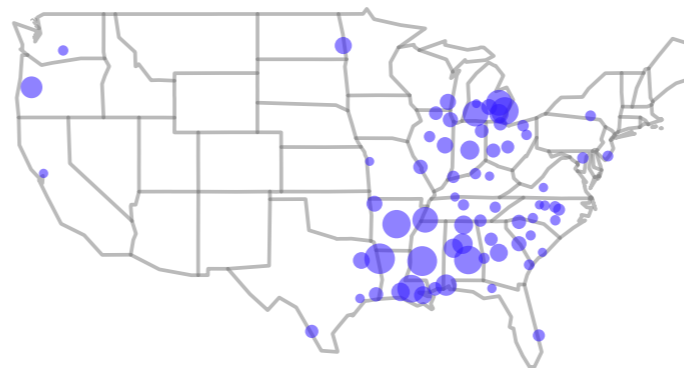
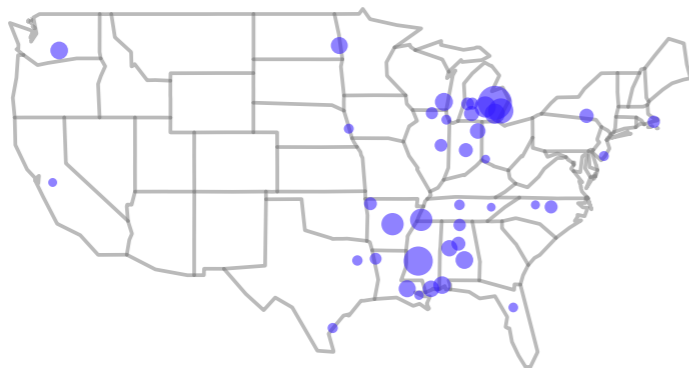
weeks 51–100

weeks 101–150

af



ikr



Do diffusion patterns follow geographic and demographic similarity?

Geolocated Twitter + U.S. Census data

7 TB data, 200 regions, 2600 words, 165 timesteps = 85M parameters

$$n_{w,r,t} \sim \text{Binom}(N_{r,t}, \sigma(\nu_w + \tau_{r,t} + \eta_{w,*,t} + \eta_{w,r,t}))$$

$$\eta_{w,t} \sim \text{Normal}(\mathbf{A}\eta_{w,t-1}, \mathbf{\Gamma})$$

$\mathbf{A}$  autoregressive coefficients (size  $R \times R$ )

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

- OK, so socially embedded language exists
- Any implications for natural language processing?

**TweetNLP:**  
**Part-of-speech tagging and word clusters**  
**for English-language Twitter**  
(available at <http://www.cs.cmu.edu/~ark/TweetNLP/>)

TweetMotif: Exploratory Search and Topic Summarization for Twitter.  
Brendan O'Connor, Michel Krieger, and David Ahn.  
*ICWSM 2010.*

Part-of-speech tagging for Twitter: Annotation, Features, and Experiments.  
Kevin Gimpel, Nathan Schneider, Brendan O'Connor, Dipanjan Das, Daniel Mills, Jacob Eisenstein, Michael Heilman, Dani Yogatama, Jeffrey Flanigan and Noah A. Smith.  
*ACL 2011.*

Improved Part-of-Speech Tagging for Online Conversational Text with Word Clusters.  
Olutobi Owoputi, Brendan O'Connor, Chris Dyer, Kevin Gimpel, Nathan Schneider and Noah A. Smith.  
*NAACL 2013.*

# NLP on social media's own terms

ikr	smh	he	asked	fir	yo	last
[Redacted]						
name	so	he	can	add	u	on
[Redacted]						
fb	lololol					
[Redacted]						

- Is this “noisy text”?
- Any NLP system, starting with POS tagging, needs different models/resources than traditional written English
  - Annotate ~2300 tweets
  - Train word clusters on 56 million tweets, use as features

# NLP on social media's own terms

ikr	smh	he	asked	fir	yo	last
!	G	O	V	P	D	A
name	so	he	can	add	u	on
N	P	O	V	V	O	P
fb	lololol					
^	!					

w fo fa fr fro ov fer **fir** whit abou aft serie fore fah fuh w/her w/that fron isn agains

“non-standard prepositions”

yeah yea nah naw yeahh nooo yeh noo noooo yea **ikr** nvm yeahhh nahh nooooo

“interjections”

facebook **fb** itunes myspace skype ebay tumblr bbm flickr aim msn netflix pandora

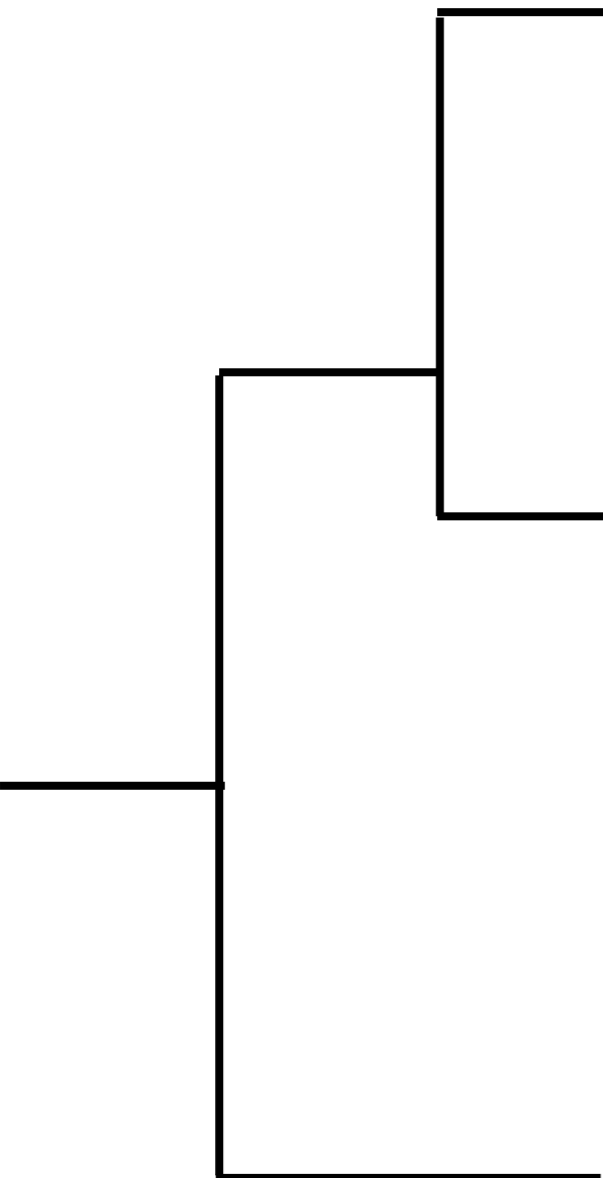
“online service names”

**smh** jk #fail #random #fact smfh #smh #winning #realtalk smdh #dead #justsaying

“hashtag-y interjections”??



- Emoticons etc.  
(Clusters/tagger useful for sentiment analysis: NRC-Canada SemEval 2013, 2014)



:d ^^ =d \*-\* :-d \o/ :dd \m/ 8d \*--\* \*\_\* u.u :ddd ;;) \*.\* o/ ;3 =))) \*---\* \('▽`)/ n\_n b-) (^\_^)  
 ^o^ :dddd ;dd \*\_\* :)))))) \*----\* d/ \o \: =dd n.n -q \*\_\* :33 :dddd :od -n \*-----\* xdddd  
 <URL-crunchyroll.com> ^^v (x \= =:) \*-----\* \0/ (~\_~")

; ) :p :-) xd ;-) ;d ( ; :3 ;p =p :-p =)) ;] xdd 😊 #gno xddd >: ) ; -p >:d ☐ 8-) ☐ 😊 ☐ ; -d ☐ 😊 [ ;  
 ☐ :^ ) =)))) ; -)) <URL-seismic.com> :pp :~) x'd :op >:p ;^ ) >:] =)))) :> ) <URL-hstl.co> ;)) )  
 ;~) toort >:3 #eden ;pp

: ) (: =) :)) :] ☺ :') =] ^\_^ :)) ^.^ [ : ;)) 😊 ((: ^\_\_^ (= ^-^ :)) 😊 👍 ☐ :-)) 😊 🙌 ^\_\_^ (: : }  
 :)) ☐ 😊 🙌 ☐ ☐ 😊 :") :]] ☐ =]] 😊 ☐ ☐ ü ;)) [= (-: ^\_\_^ ;') :-)) (((:

:o o\_o « o.o xp ;o .\_. t.t t\_t #wtf #lol o: x\_x =o 0\_o dx o\_0 :-o ~-~ --" 0\_0 o\_o »» u\_u #help  
 --' =3 (-\_-) -) #confused ☐ #omg ~-~ t^t otl #igetthatalot 🤪 xdddd o\_\_o @@ cx t\_\_t d8 ☐  
 :{ t\_\_t ----- #whodoesthat e\_e :oo

:( :/ -\_- -.- :- ( :'( d: :l :s -\_- =( =/ >.< -\_- - :-/ </3 :\ -\_- - ;( /: ☹ :(( >\_< =[ :[ #fml 😞 -  
 \_\_\_\_\_ - =\ >:( 😞 -,- >> >:o ;/ 🤪 d; .- - \_\_\_\_\_ - >\_> :((( -\_- " =s ☐ ;\_ ; #ugh :-\ =.= ☐ -  
 \_\_\_\_\_ -

x xx xxx xxxx xxxxx qt xox xxxxxx xxxxxxx xxxxxxxx #pawpawty xxxxxxxxxx xxxxxxxxxx  
 #1dfamily #frys #1dqa xxxxxxxxxxxx #askliam #dcth xxxxxxxxxxxxxx #askniiall \*rt  
 #jbinpoland xxxxxxxxxxxxxx #askharry x-x #wiimoms xxxxxxxxxxxxxxxxxxxx oox #wlf #nipclub  
 +) 1dhq xxxxxxxxxxxxxxxxxxxx #20peopleilove <URL-paidmodels.com> yart #jedreply  
 #elevenestime <URL-shrtn.us> #askzayn xxxxxxxxxxxxxxxxxxxx #wineparty +9  
 #amwritingparty #tweepletuesday #soumanodomano <URL-today.com> #twfanfriday 22h22

<3 ♥ xoxo <33 xo <333 ♥ ♥ #love s2 <URL-twition.com> #neversaynever <3333 #swag  
 x3 #believe #100factsaboutme ♥♥ 🤪 <3<3 <33333 #blessed xoxoxo 😊 #muchlove  
 #salute xoxox ♥♥♥ #excited 🌟 ☐ #happy #leggo #cantwait <3<3<3 #loveit <333333  
 #please #dailytweet #thanks 🙏 (~\_~) 💜 #yay #thankyou #loveyou {} ε~) #nsn #iloveyou

# Subject-AuxVerb constructs

[Contraction  
splitting?]

[Mixed]

i'd you'd we'd he'd they'd she'd who'd i'd u'd youd you'd iwould theyd  
icould we'd i`d #whydopeople he'd i´d #iusedto they'd i'ld she'd  
#iwantsomeonewhowill i'de imust a:i'd you`d yu'd icud l'd

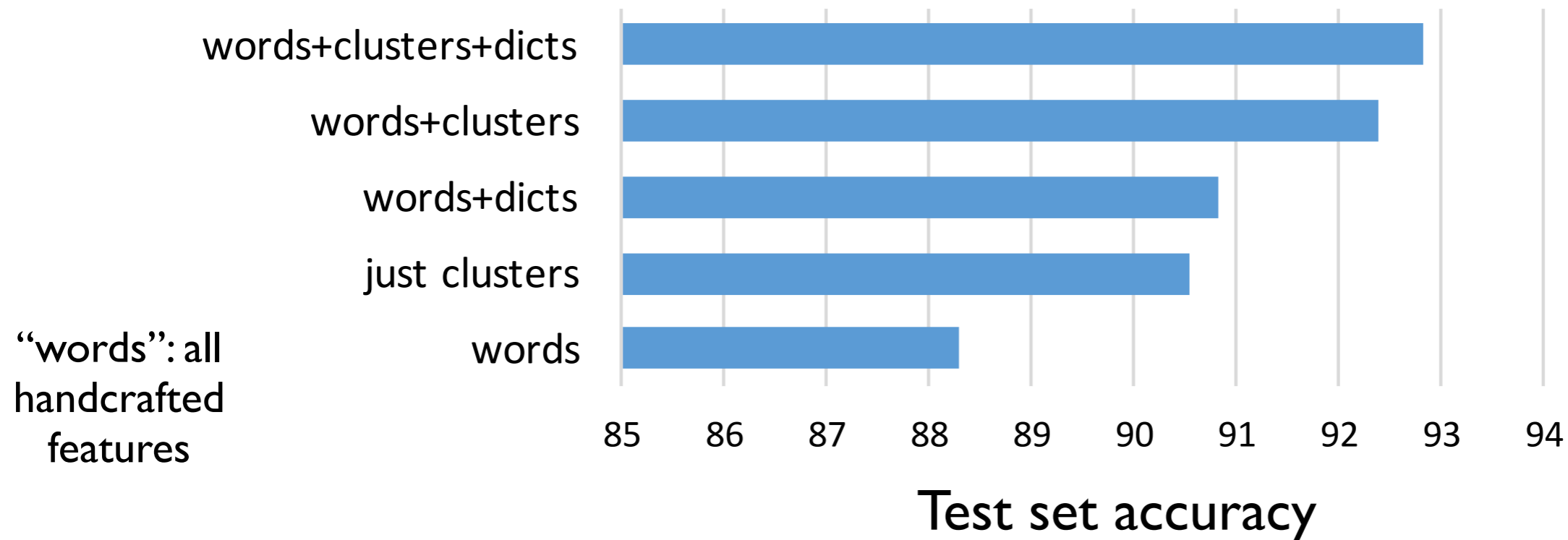
ill ima imma i'ma i'mma ican iwanna umma imaa #imthetypeto iwill  
amma #menshouldnever igotta #whywouldyou #iwishicould  
#sometimesyouhaveto #thoushallnot #ihatewhenpeople illl  
#thingspeopleshouldnotdo #howdareyou #thingsgirlswantboystodo  
im'a #womenshouldnever #thingsblackgirlsdo immma iima  
#ireallyhatewhenpeople ishould #thingspeopleshouldntdo #irefusetto itl  
#howtospoilahoodrat iwont imight #thingsweusedtodoaskids ineeda  
#thingswhitepeopledo we'l #whycantyoutjust #whydogirls  
#everymanshouldknowhowto #ushouldnt #howtopissyourgirloff  
#amanshouldnot #uwannaimpressme #realfriendsdont immaa  
#ilovewhenyou

you'll we'll it'll he'll they'll she'll it'd that'll u'll that'd youll ull you'll itll  
there'll we'll itd there'd theyll this'll thatd thatll they'll didja he'll it'll  
yu'll she'll youl you`ll you'l you´ll yull u'l it'l we´ll we`ll didya that'll  
it'd he'l shit'll they'l theyl she'l everything'll he`ll things'll u'll this'd

i'll i´ll i'l i`ll i´ll i'lll l'll i`ll i"ll -i'll /must @pretweeting she`ll



# Clusters help POS tagging



- A little annotation + lots of data
- Unsupervised word representation learning (clusters, embeddings) is a crucial technique in NLP

- **Where do nonstandard terms come from?**

imma



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

imma



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



**r.i.p spikeeee** @xEnvyme\_alex · 2m

Feel like **imma** have it up way more my sophomore year than I did freshman year



**Angelica** @BrowncoatAnge · 2m

**Imma** start unfollowing all these celebs crying about a Lion and not even remotely interested in the struggle of the Black community.



**Billi3 J3an** @Misskooki3 · 2m

Niggas thought they wanted to pimp me before they noticed bitch **imma** boss.. I pimp the nigga who brought u here



**kye** @hippys0ul · 2m

K **imma** just scroll my tweets to the day abel followed. 🤔❤️

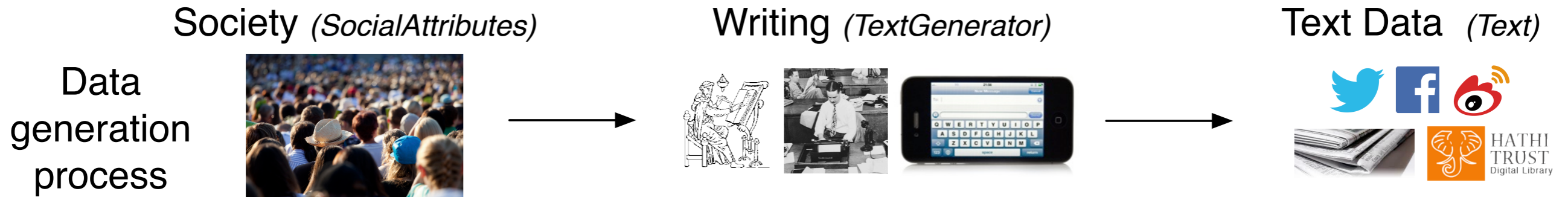


**LifeΣtyleTr3** @LifeStyleTr3 · 2m

@Charismatic\_Cee ay shoot me those pictures you calm you got of us n **imma** post em on the gram !



# TextGenerator(SocialAttributes) → Text



## What to analyze:

Social phenomena in social media datasets

- Political speech under Chinese censorship
- Events in international relations
- Social factors in language use

## How to analyze:

NLP capabilities we need to do these better

- Part of speech tagging
- Entity extraction
- Syntactic, semantic parsing

**What social bias exists in NLP models?**

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

Su Lin Blodgett



Lisa Green



Brendan O'Connor



*EMNLP 2016*

**What social bias exists in NLP models?**

# Dialect

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

RETWEETS

3

LIKES

42



1:08 AM - 8 Jul 2016



# Dialect

SAE:  
*he is woke af*



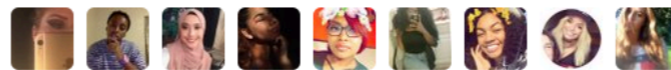
he woke af smart af educated af daddy af  
coconut oil using af GOALS AF & shares food af

RETWEETS

3

LIKES

42



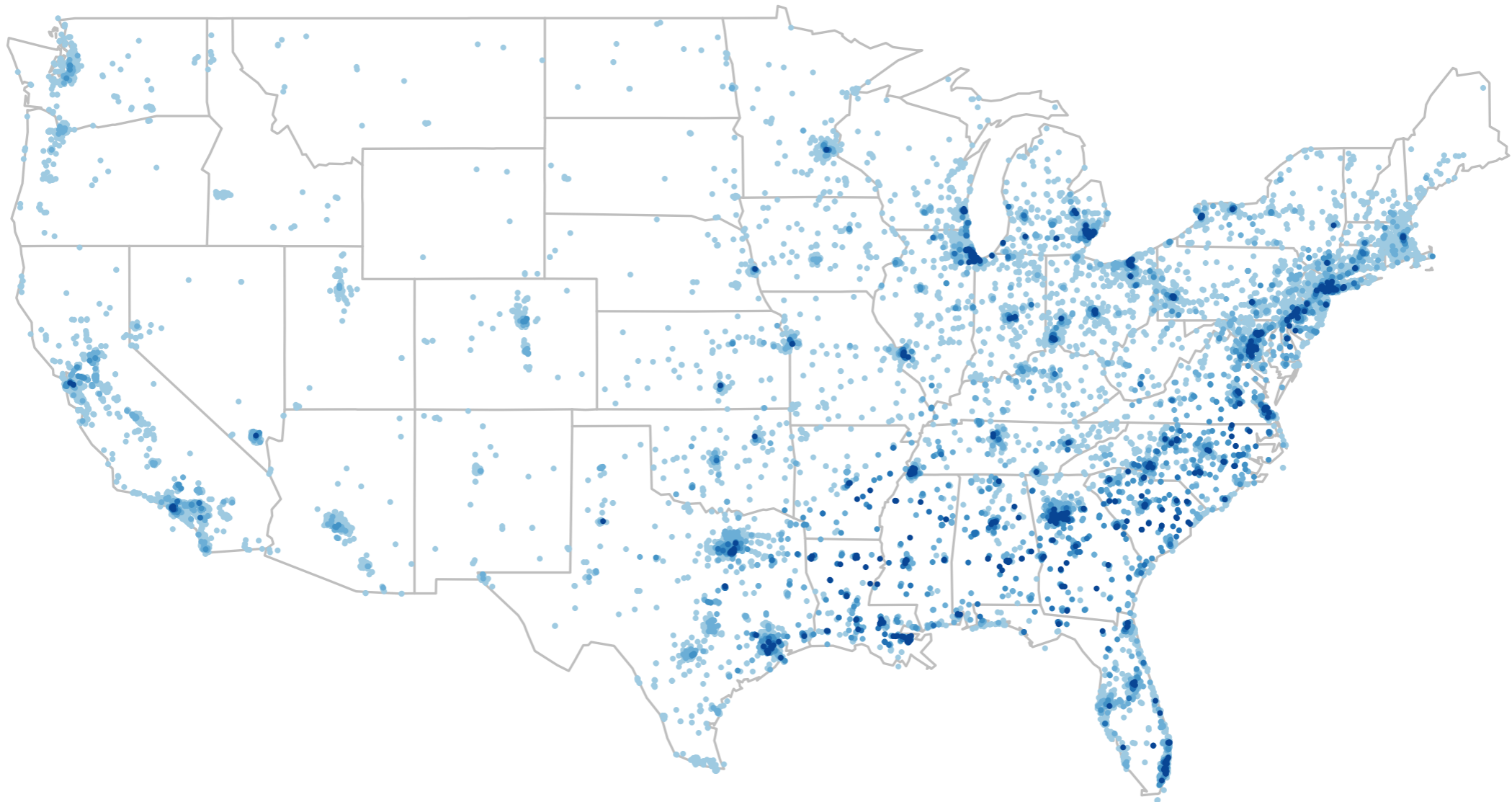
1:08 AM - 8 Jul 2016





# Why is social media different?

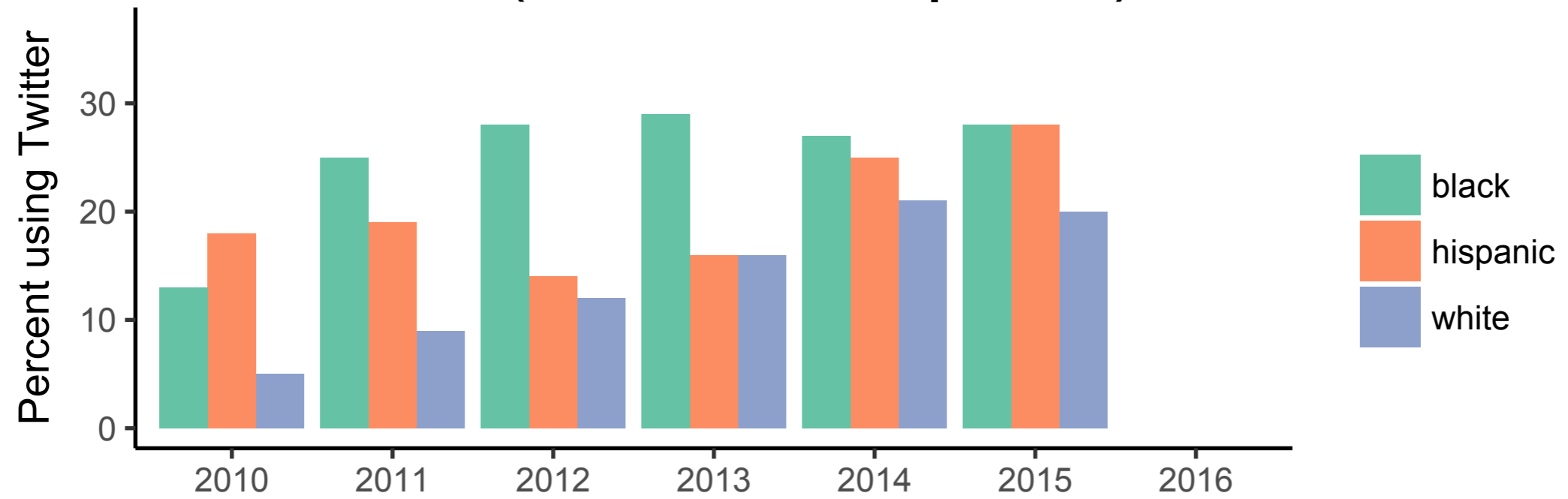
- Internet speech?
- Pre-existing dialectal English?
  - Geographic patterns of word usage often reveal relationships to race, ethnicity etc.
  - African-American English in Twitter  
*[Eisenstein 2013, Jorgensen et al. 2015, Jones 2015]*



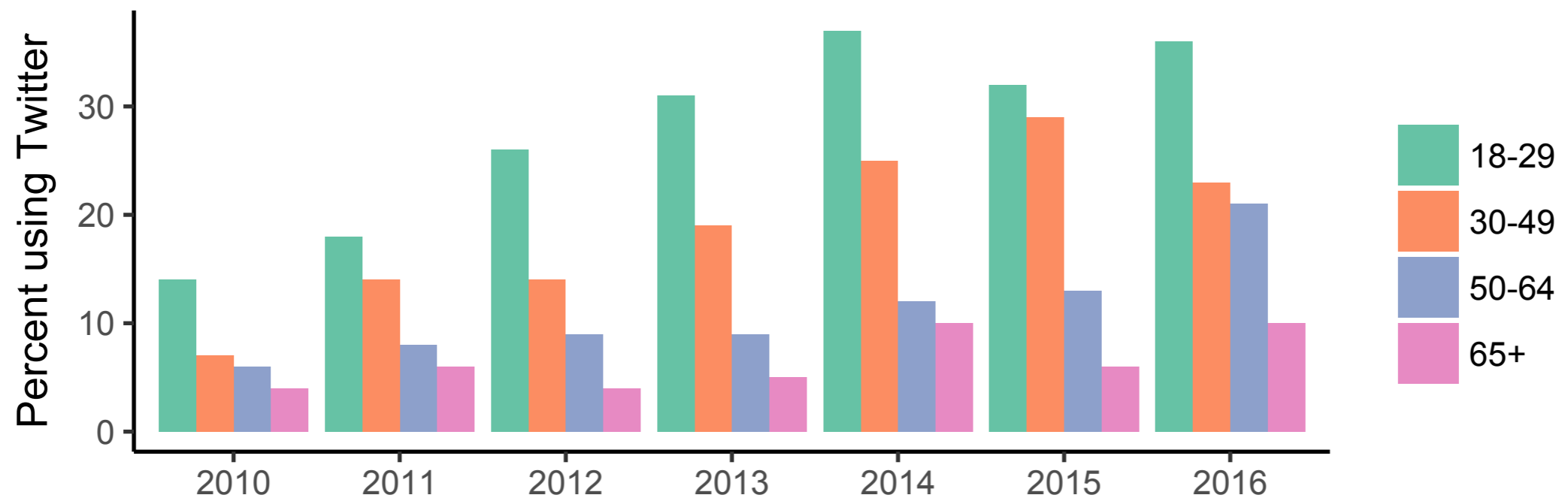
# Youth, minorities on Twitter

[Pew Research]

## P(use twitter | race)



## P(use twitter | age)



- From U.S. Census data and geo-located tweets: identify demographic-specific terms and messages via probabilistic model
- Validate African-American-associated corpus against linguistics literature on African-American English
- Investigate racial disparities in natural language processing tools

# Associating geolocated tweets with demographics



**King Me**  
@tmac\_datboi

 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!!!

# Associating geolocated tweets with demographics



**King Me**  
@tmac\_datboi



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

block group 010730039001

Bored af den my phone finna die!!!

block group 010730058003

# Associating geolocated tweets with demographics

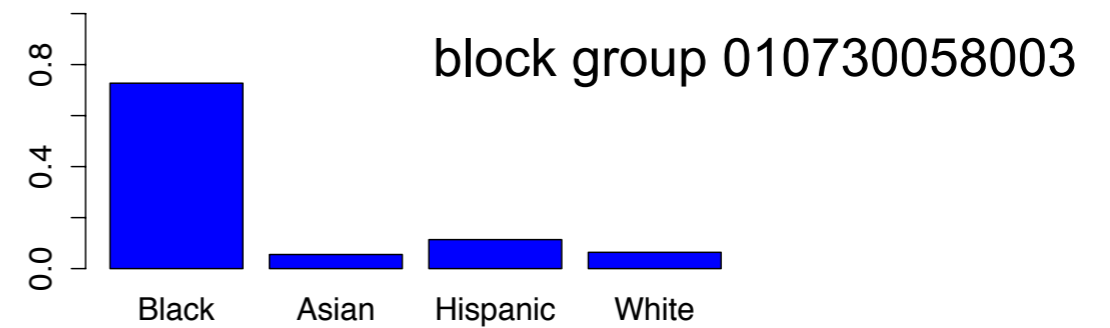
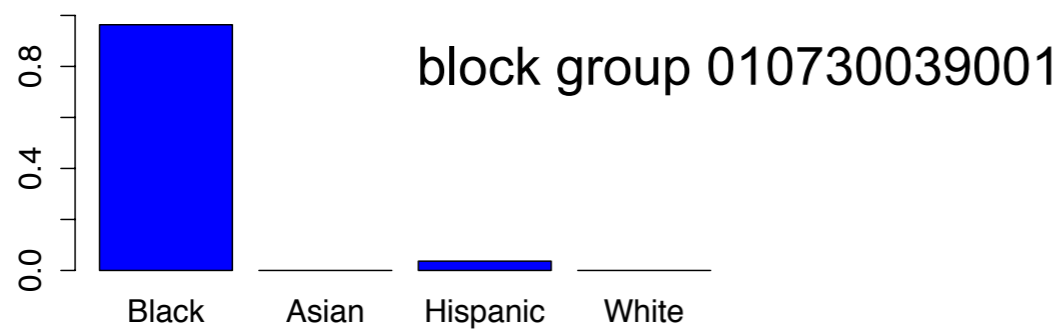


**King Me**  
@tmac\_datboi

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!!!



# Associating geolocated tweets with demographics

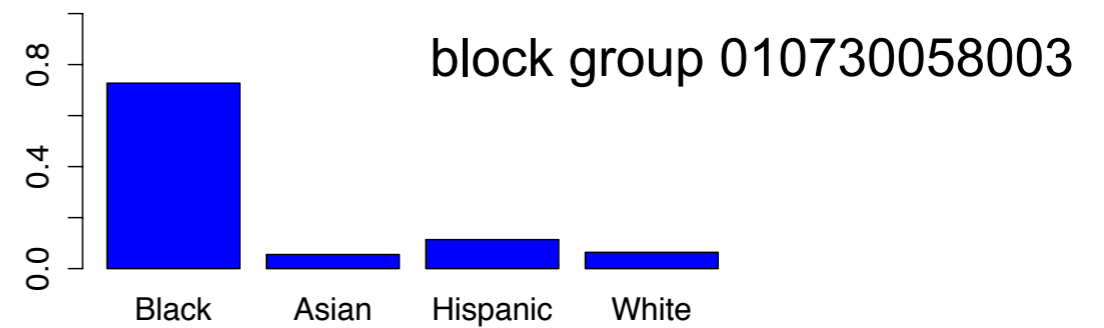
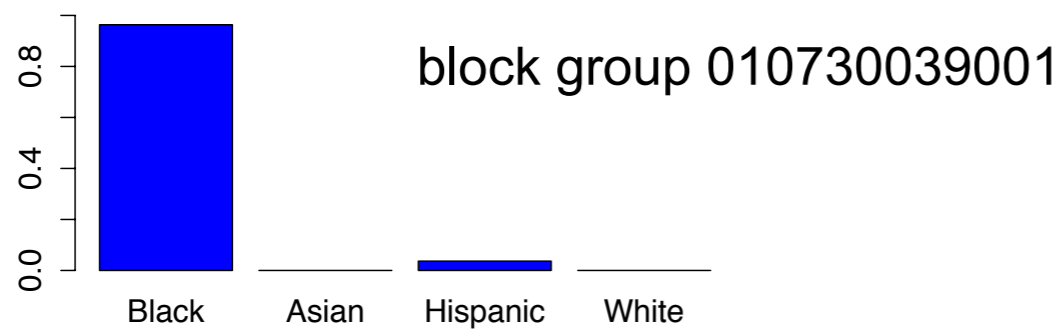


King Me  
@tmac\_datboi

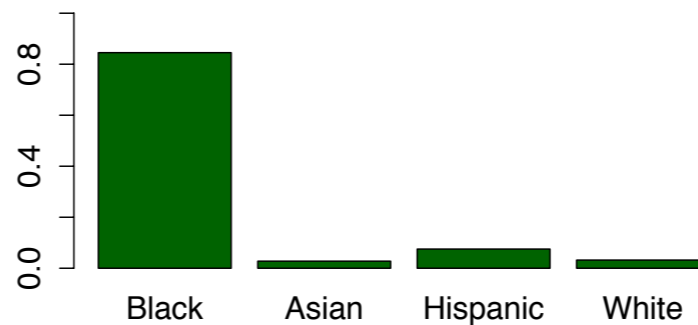


he woke af smart af educated af daddy af  
coconut oil using af GOALS AF & shares food af

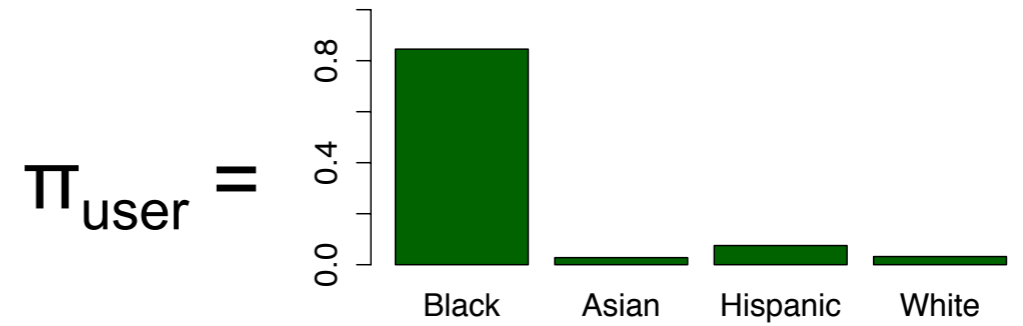
Bored af den my phone finna die!!!



$$\pi_{\text{user}} =$$



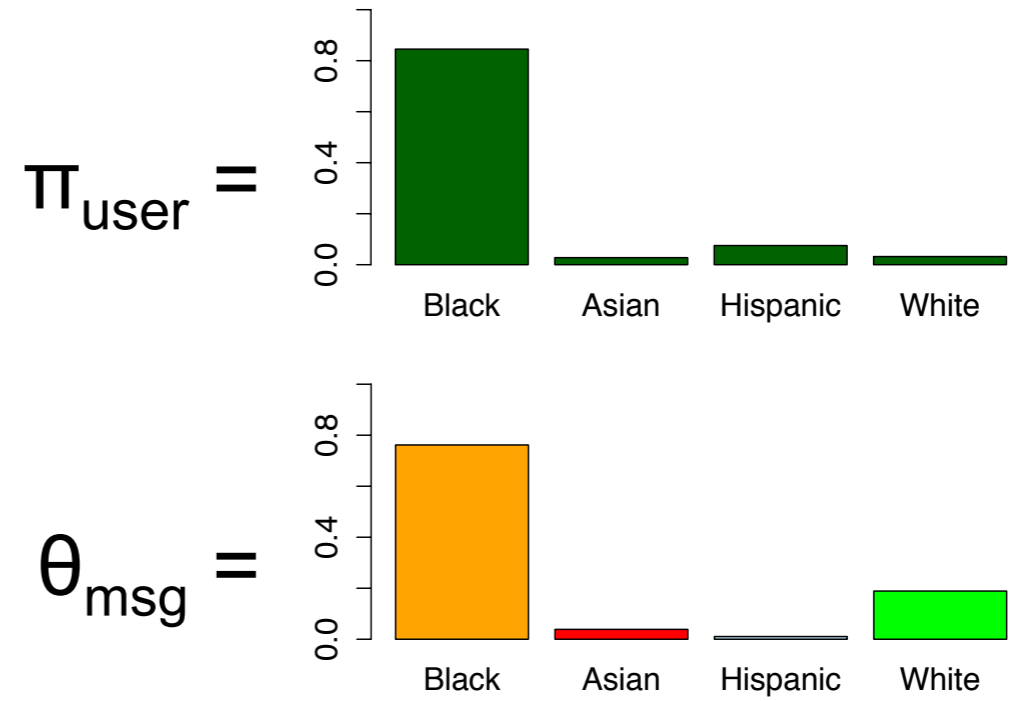
# Mixed membership model





# Mixed membership model

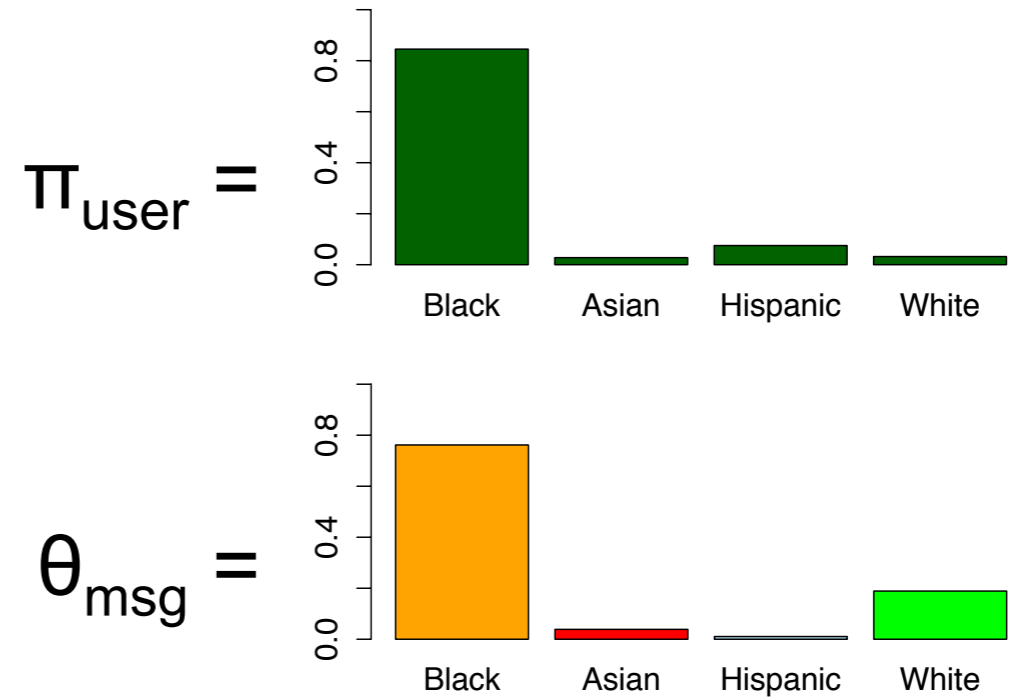
$$\theta_{\text{msg}} \sim \text{Dir}(\alpha\pi)$$



# Mixed membership model

$$\theta_{\text{msg}} \sim \text{Dir}(\alpha\pi), \quad z \sim \theta_{\text{msg}}, \quad w \sim \phi_z$$

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

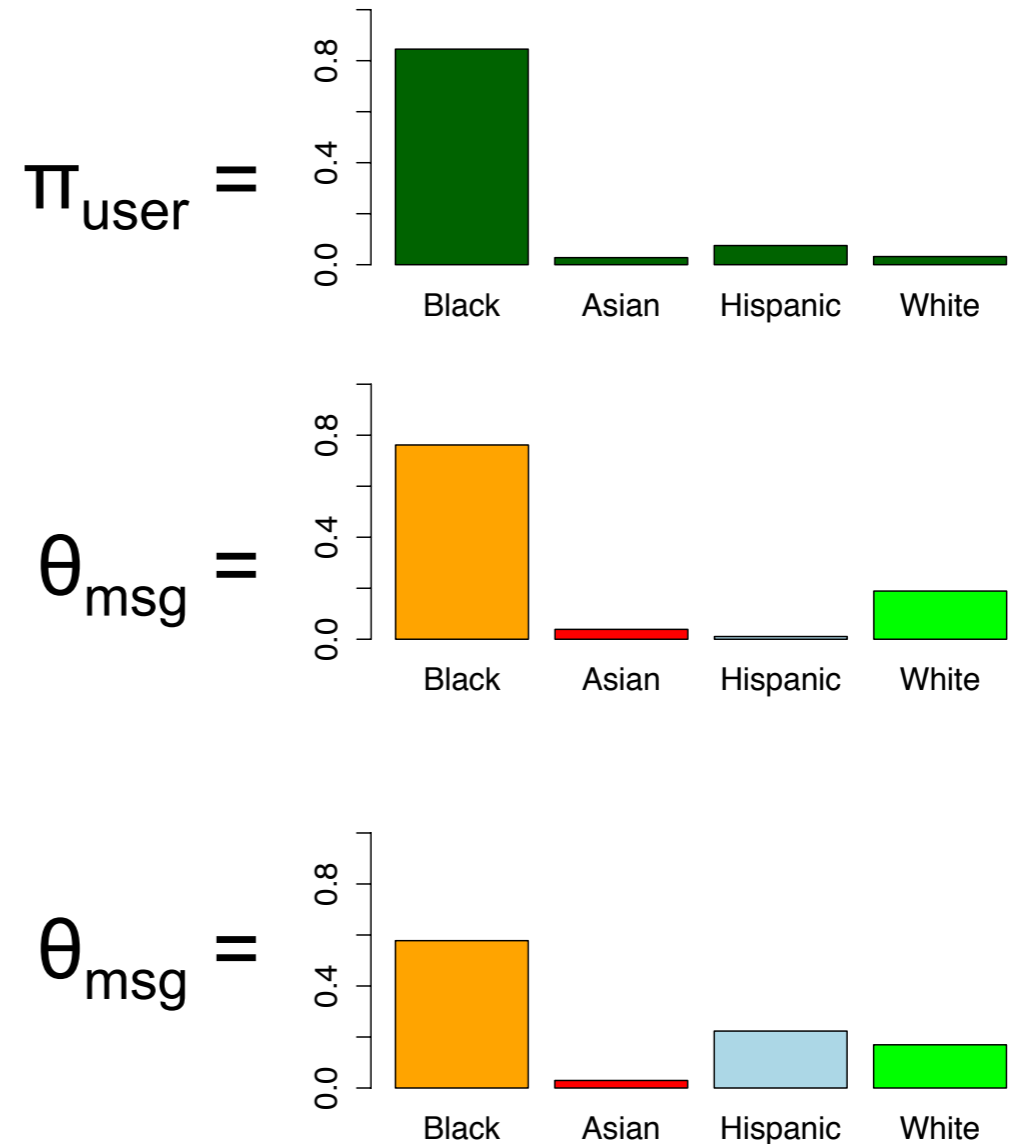


# Mixed membership model

$$\theta_{\text{msg}} \sim \text{Dir}(\alpha\pi), \quad z \sim \theta_{\text{msg}}, \quad w \sim \phi_z$$

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!!!



# Mixed membership model

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

$m_1$

Bored af den my phone finna die!!!

$m_2$

# Mixed membership model

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

$m_1$

Bored af den my phone finna die!!!

$m_2$

Word	AA	Asian	Hisp.	White
woke	1	0	0	0
af	6	0	0	0
educated	0	0	0	1
...				

# Mixed membership model

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

$m_1$

Bored af den my phone finna die!!!

$m_2$

Word	AA	Asian	Hisp.	White
woke	1	0	0	0
af	6	0	0	0
educated	0	0	0	1
...				

Message	AA	Asian	Hisp.	White
$m_1$	7	0	0	2
$m_2$	2	0	1	1

# Mixed membership model

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

$m_1$

Word	AA	Asian	Hisp.	White
woke	1	0	0	0
af	6	0	0	0
educated	0	0	0	1
...				

Bored af den my phone finna die!!!

$m_2$

Message	AA	Asian	Hisp.	White
$m_1$	7	0	0	2
$m_2$	2	0	1	1

User	AA	Asian	Hisp.	White
$u_1$	9	0	1	3

# Corpus creation and linguistic validation

- Beyond unigrams: creation of user-level topic-aligned corpora



# 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  $r_{AA}(w) \geq 2$ ,  $r_{white}(w) \geq 2$ : AA- and white-aligned words

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

# Phonological analysis

- Calculate  $r_{AA}(w)$  for 31 phonological variants illustrated through nonstandard spellings

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

# Phonological analysis

- Calculate  $r_{AA}(w)$  for 31 phonological variants illustrated through nonstandard spellings
- For 30/31 variants:  $r \geq 1$

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

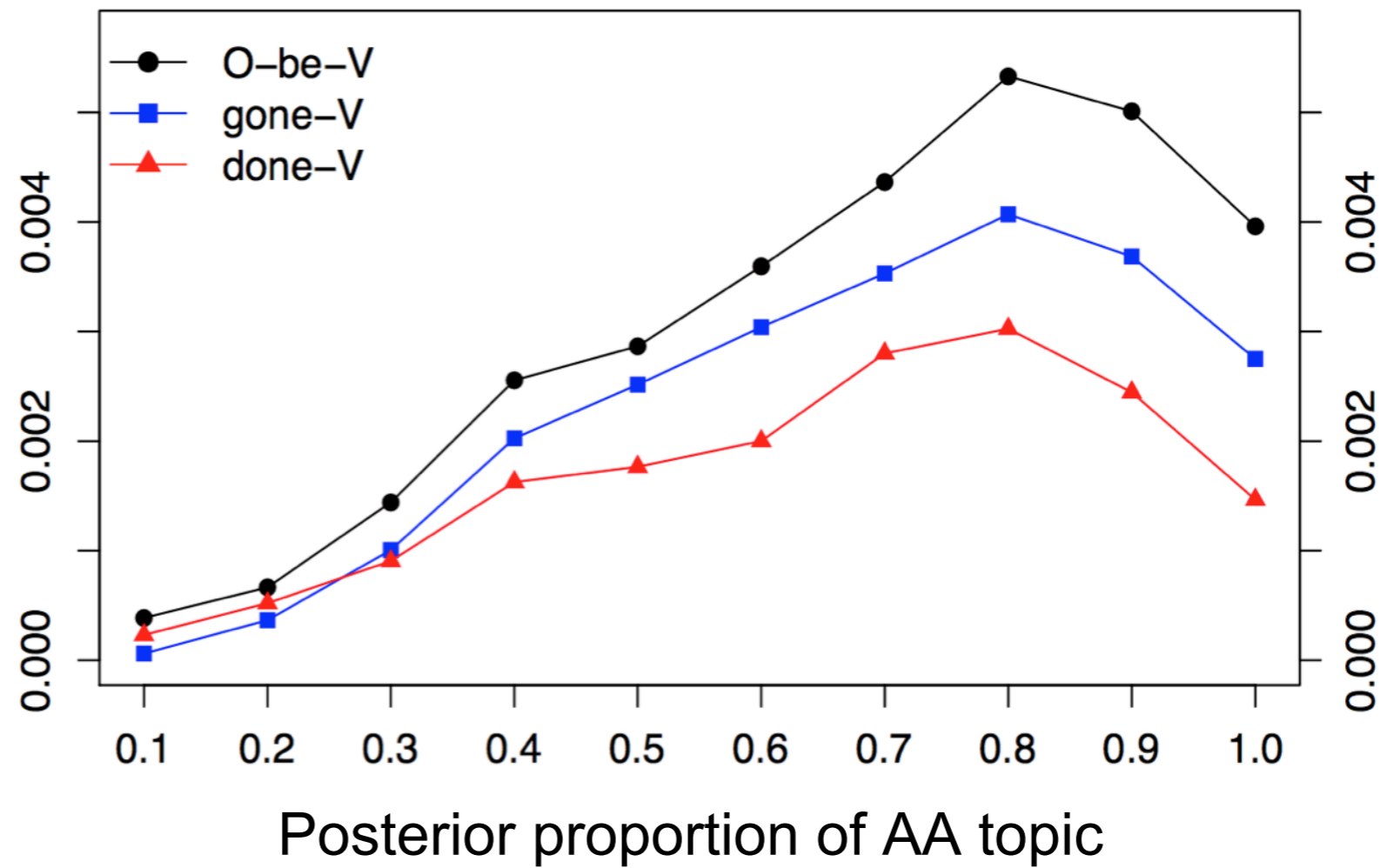
# Syntactic analysis

- Select 3 well-known AAE verbal markers
- Search for sequences of unigrams and POS tags

<b>Construction</b>	<b>Example</b>
<i>O-be/b-V</i>	<i>I be tripping bruh</i>
<i>gone/gne/gon-V</i>	<i>Then she gon be single Af</i>
<i>done/dne-V</i>	<i>I done laughed so hard that I'm weak</i>

# Syntactic analysis

Proportion of tweets with construction



# Historical vs. Online?

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**

---



# Historical vs. Online?

## 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

$$\frac{P(\text{dat} \mid AA)}{P(\text{dat} \mid \neg AA)} = 5.9$$

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

# Historical vs. Online?

## 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

$$\frac{P(\text{dat} \mid AA)}{P(\text{dat} \mid \neg AA)} = 5.9$$

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

## POS taggers: standard vs. designed for Twitter

---

CoreNLP	dey/ <b>NN</b> (PRP) b/ <b>NN</b> (VBZ) ' /Punct longs/ <b>NNS</b> (VBZ) to/TO dat/ <b>VB</b> (DT) gent/JJ ' /Punct man/NN ahaid/ <b>VBN</b> (RB)
---------	--

---

ARK	dey/Pro b'longs/Verb to/Prep <b>dat</b> /Det gent'man/Noun ahaid/Adv
-----	--

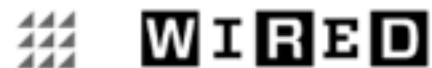
---

# Examining NLP tools - dependency parsing

- Compare annotated parses to systems' output parses

# Examining NLP tools - dependency parsing

- Compare annotated parses to systems' output parses



Google Has Open Sourced SyntaxNet, Its AI for Understanding Language

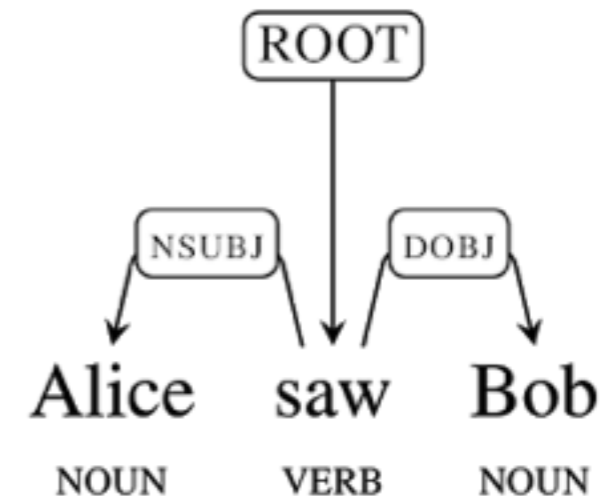


Google Research Blog

[Announcing SyntaxNet: The World's Most Accurate Parser Goes Open Source](#)

Thursday, May 12, 2016

Posted by Slav Petrov, Senior Staff Research Scientist



## 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)

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

## 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.

# Examining NLP tools - language identification

- Language identification - key step in NLP pipelines

# 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'''\n>>> langid.classify(s)\n('da', 0.99999999993212958)
```

```
>>> s = 'Bored af den my phone finna die!!!'\n>>> langid.classify(s)\n('da', 0.9999968001354156)
```



# Examining NLP tools - language identification

- Language identification - key step in NLP pipelines

	<b>AA-Aligned</b>	<b>White-Aligned</b>
<i>langid.py</i>	13.2%	7.6%

Proportion of messages classified as  
non-English

# Examining NLP tools - language identification

- Language identification - key step in NLP pipelines

	<b>AA-Aligned</b>	<b>White-Aligned</b>
<i>langid.py</i>	13.2%	7.6%
<b>Twitter</b>	24.4%	17.6%

Proportion of messages classified as  
non-English

# Examining NLP tools - language identification

- Solution: build ensemble classifier to augment langid.py
- Given a message, classifier:
  - Calculates langid.py's prediction
  - If prediction is English, return English
  - If not English, return English if our model's (AA + white + Hispanic) posterior probabilities  $\geq 0.9$
  - Otherwise, return langid.py's prediction

# Examining NLP tools - language identification

- Solution: build ensemble classifier to augment langid.py

<b>Message set</b>	<i>langid.py</i>	<b>Ensemble</b>
High AA	80.1%	99.5%
High White	96.8%	99.9%
<i>General</i>	88.0%	93.4%

Imputed recall of English messages for  
2014 messages

- Develop a model leveraging demographic correlations to generate dialectal corpora
- Corpus reproduces well-known dialectal phenomena
- Demonstrate disparity in performance by two kinds of NLP tools
- Provide ensemble classifier augmenting existing tools with our model

# 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