

We some young kings: Communities, age, and African American English online

Ian Stewart, Dartmouth College
ian.b.stewart.14@alum.dartmouth.org
Advised by Sravana Reddy, James Stanford

Motivation

The structure of AAE's syntax is well understood¹ but its demographic distribution remains debated. This study seeks to:

1. Develop a method to reliably extract short-range syntactic constructions from Twitter data via regular expressions.
2. Corroborate geographic² and/or demographic variation in syntax usage.

Data Collection

Geotagged Tweets (mined July - December 2013)

Filtering (users > 20 tweets)

228 million tweets, 1 million users

Content

Metadata

POS Tags³

Regular Expressions

ZIP Codes
(linked to US Census demographic data)

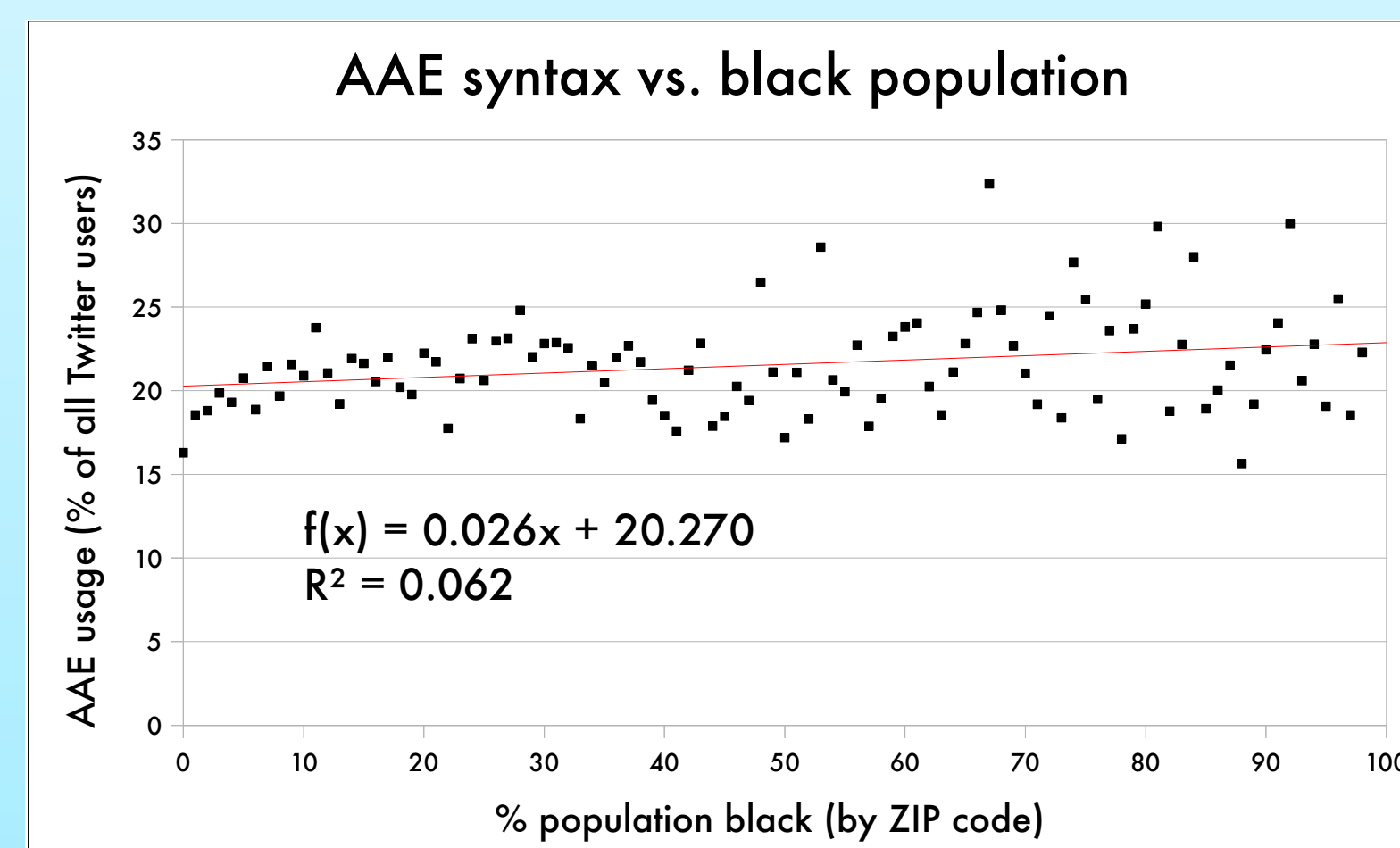
Gender
(estimated using 1995 Social Security Administration name database)

Construction ¹	Example	Frequency
ass camouflage construction (ACC)	<i>divorced her ass</i>	578,422
continuative steady	<i>steady washing her face</i>	15,979
copula deletion	<i>you mad</i>	1,545,024
future finna	<i>I'm not finna stress myself</i>	338,143
habitual be	<i>he be drinking that beer</i>	588,752
negative concord	<i>I ain't got no time</i>	455,531
negative inversion	<i>ain't nobody scared of you</i>	97,042
null genitive	<i>everybody & they mama</i>	296,358
past complete done	<i>Hunter done killed her</i>	87,346
remote past been	<i>Chris been told me</i>	16,998

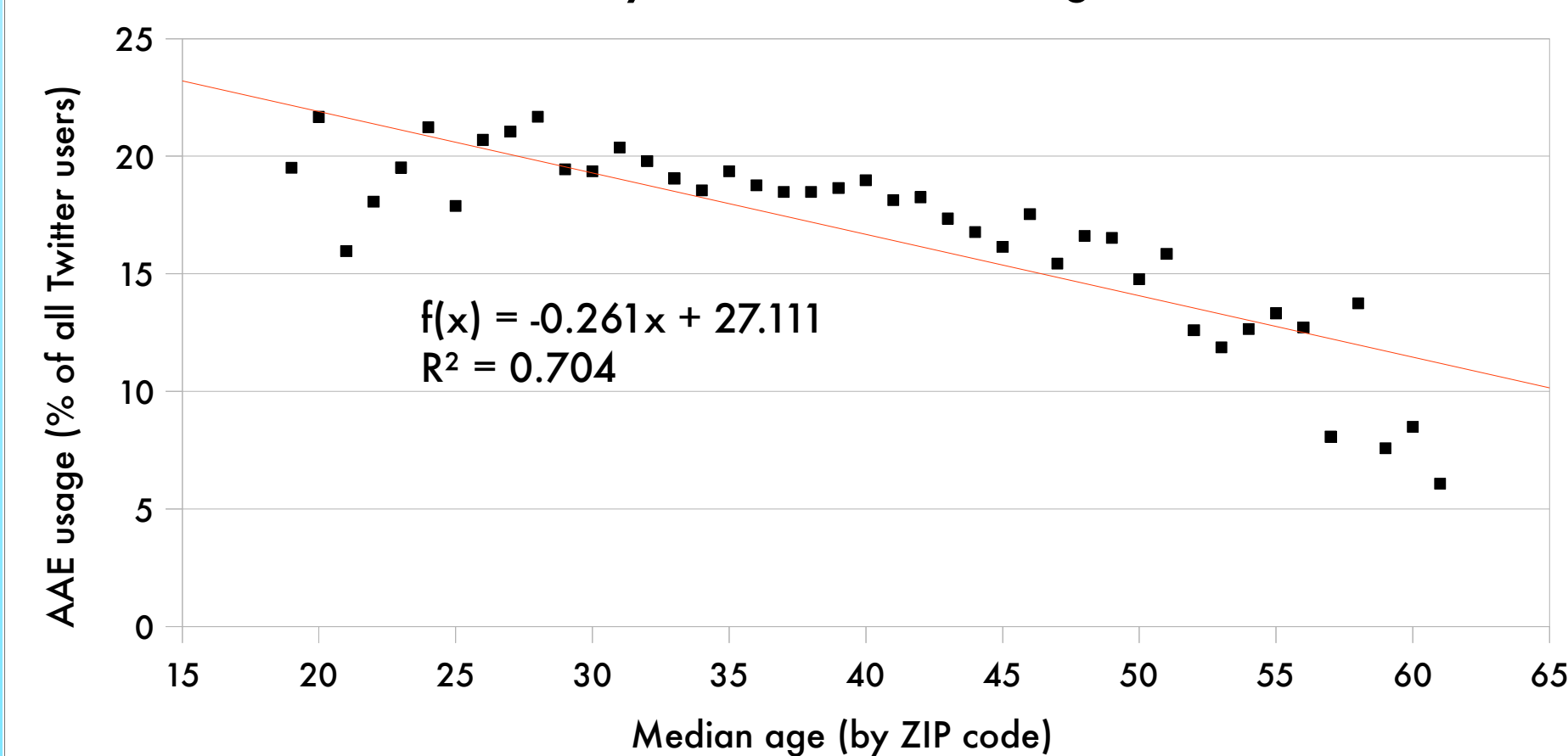
Results

Ethnicity?

- Weak positive correlation between black population and AAE syntax usage.
- Unreliable trends for other ethnic demographics and geography (ex. by state).



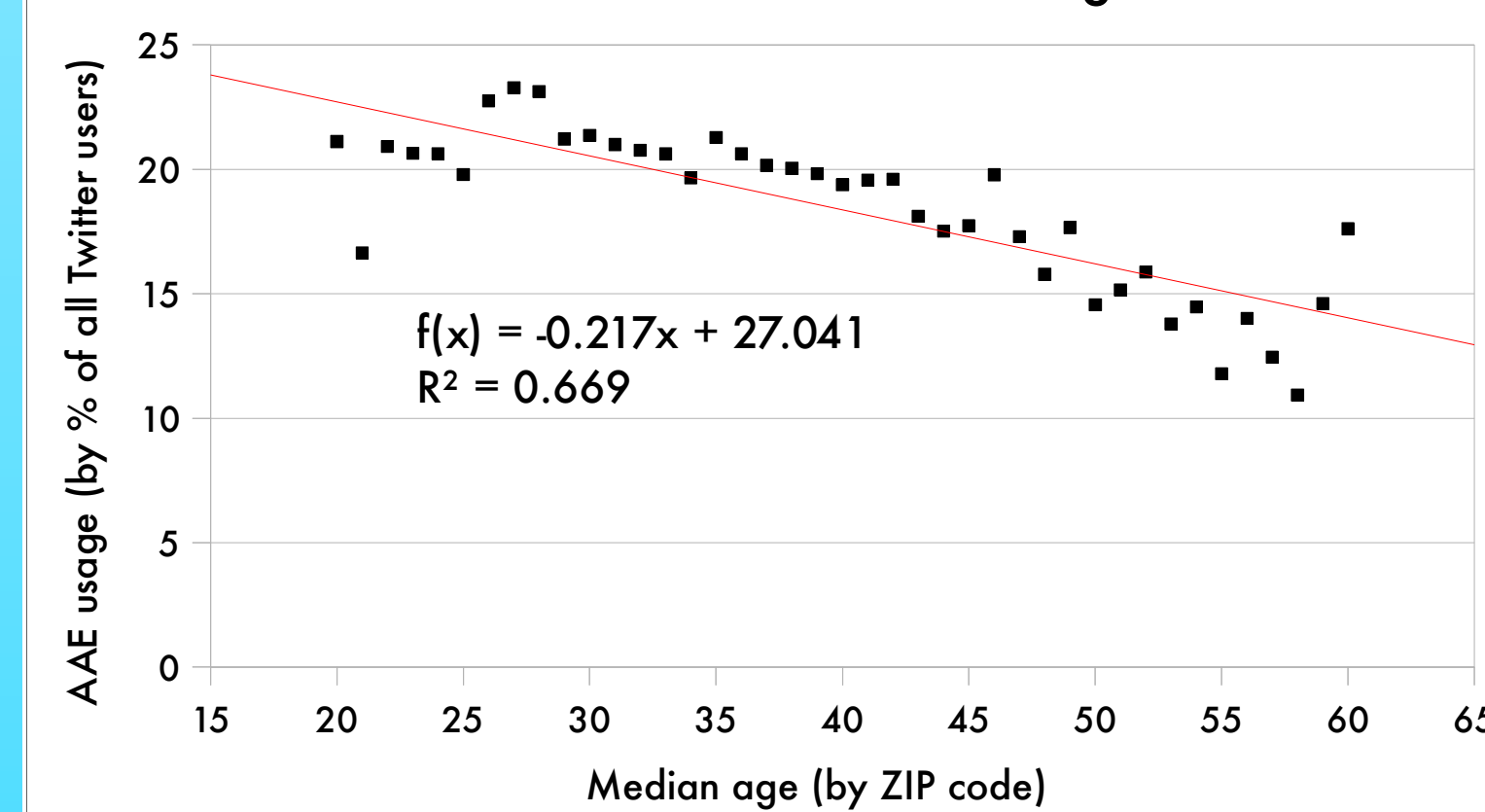
AAE syntax vs. median age



Median Age!

- Clear negative correlation, especially from ages 20-60 (>95% of data).
- Growth in usage through 20s followed by steady decline.

AAE lexemes vs. median age



Confirmation.

- Same test performed for AAE lexemes⁴ (ex. <talmbout>) and baseline lexemes (ex. <the>).
- Identical negative trend in AAE lexemes and absence in baseline lexemes.

Conclusions and Future Directions

- No significant variation in AAE syntax usage with respect to ethnic demographics.
 - ▶ Potential cause: African American users' overrepresentation in Twitter⁵.
- Strongest trend? Median age.
 - ▶ Community-level curvilinear principle, even outside of towns with expected Twitter representation.
 - ▶ Age might be correlated with social structure (college town vs. family suburb).
- Caveat: hard to separate legitimate language use from quotations.
- Next step: check for connections between median age and other community patterns.
- Determine importance of social networks versus community statistics.
 - ▶ Twitter users often group by ethnicity⁶.
 - ▶ Denser social networks in younger communities?
- Improve syntax detection accuracy.
 - ▶ 16.6% false positive rate, incalculable true/false negatives.
 - ▶ More rigorous regexes, possibly with crowdsourcing.
 - ▶ Training word-vector model for semantics (ex. ACC).
 - ▶ Bootstrap syntax from lexicon (ex. train on <talmbout> tweets).

References

1. Rickford, John R. *African American Vernacular English*. Malden, MA: Blackwell Publishers, 1999.
2. Wolfram, Walt. "Urban African American Vernacular English: morphology and syntax." *A handbook of varieties of English*. Ed. Bernard Kortmann. Walter de Gruyter, 2004.
3. Olutobi Owoputi et al. "Improved Part-of-Speech Tagging for Online Conversational Text with Word Clusters." *Proceedings of NAACL-HLT 2013*. Atlanta, GA: 2013.
4. Ofori-Atta, Akoto. "Here's Your Black Twitter Welcome Manual" *The Root*, Jan. 3 2014.
5. Duggan, Maeve and Joanna Brenner. "The Demographics of Social Media Users - 2012." Washington, D.C.: Pew Research Center, 2013.
6. Shane Bergsma et al. "Broadly Improving User Classification via Communication-Based Name and Location Clustering on Twitter." *Proceedings of NAACL-HLT 2013*, Atlanta, GA, 2013.