

Tweetment Effects on the Tweeted: An Experiment to Reduce Twitter Harassment

Kevin Munger

NYU

December 1, 2016

Prepared for Yale Human Nature Lab

Project Outline

- Find harassers on Twitter

Project Outline

- Find harassers on Twitter
- Randomize assignment 2x2
 - ▶ In-group v Out-group (race of bot)
 - ▶ High Status v Low Status (number of followers)

Project Outline

- Find harassers on Twitter
- Randomize assignment 2x2
 - ▶ In-group v Out-group (race of bot)
 - ▶ High Status v Low Status (number of followers)
- Offensiveness as outcome variable

Project Outline

- Find harassers on Twitter
- Randomize assignment 2x2
 - ▶ In-group v Out-group (race of bot)
 - ▶ High Status v Low Status (number of followers)
- Offensiveness as outcome variable
- See how long the treatment persists

Why does it happen?

- Prejudice

Why does it happen?

- Prejudice
- Social Identity model of De-Individuation (SIDE)

Why does it happen?

- Prejudice
- Social Identity model of De-Individuation (SIDE)
- Girard's theory of mimetic desire

State of the art

- Experiments in the lab
 - ▶ Convenience samples
 - ▶ Short time frame
 - ▶ In the lab

State of the art **My Approach**

- ~~Experiments in the lab~~ **Experiment in the “field”**
 - ▶ ~~Convenience samples~~ **Sample of real, consistent harassers**
 - ▶ ~~Short time frame~~ **Continuous and unbounded time frame**
 - ▶ ~~In the lab~~ **In the same context as the harassment**

Find harassers

- Needs to be fast, and accurate

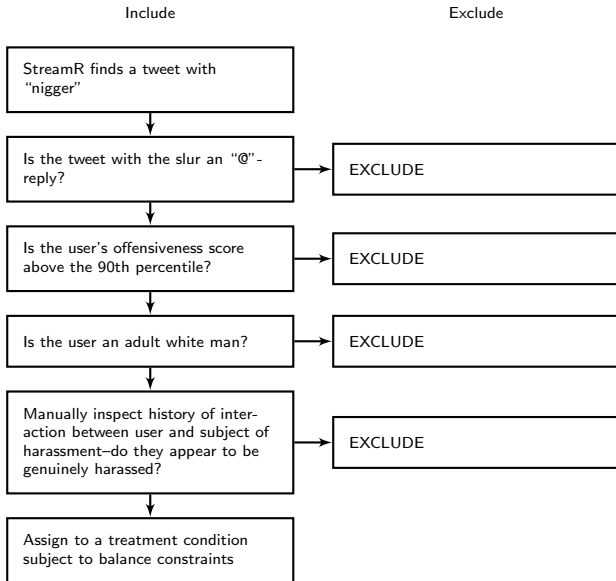
Find harassers

- Needs to be fast, and accurate
- Don't care about recall

Find harassers

- Needs to be fast, and accurate
- Don't care about recall
- In the presence of strongly offensive language, a dictionary of slurs is best (Chen et al, 2012)

Detection “Algorithm”



Apply Treatment



██████████ · 13 Sep 2015

@██████████ don't be a nigger



Rasheed ██████████

@Rasheed ██████████

@██████████ Hey man, just remember that there are real people who are hurt when you harass them with that kind of language

Treatment uptake

Greg [redacted]
@Greg [redacted]
New York, NY

TWEETS 70 FOLLOWING 39 FOLLOWERS 2 Edit profile

Tweets Tweets & replies

Greg [redacted] Retweeted
SportsCenter @SportsCenter · 11m
Michael Oher's "Blind Side" family joined him on the field to celebrate his team advancing to the Super Bowl. es.pr/1QwVGnw
Retweet 510 Like 895 ... View summary

Greg [redacted] @Greg [redacted] 15 Sep 2015
@ [redacted] Hey man, just remember that there are real people who are hurt when you harass them with that kind of language
Reply Retweet Like ... View conversation

Who to follow · Refresh · View all

- NYPD NEWS** @NYPDnews · Followed by NYC Mayor's O...
Follow
- Adam Scheffer** @Adam...
Follow

Find friends

Trends Change

Hypotheses

All hypotheses have been pre-registered through EGAP.

Hypothesis

The ranking of the magnitudes of the decrease in harassment will be:
In-group/High status > $\frac{\text{In-group/Low status}}{\text{Out-group/High status}}$ > *Out-group/Low status.*

Hypotheses

All hypotheses have been pre-registered through EGAP.

Hypothesis

*The ranking of the magnitudes of the decrease in harassment will be:
In-group/High status > $\frac{\text{In-group/Low status}}{\text{Out-group/High status}}$ > Out-group/Low status.*

Table: Experimental Design and Hypothesized Effect Sizes

	In-group	Out-group
Low followers	Medium effect	Small effect
High followers	Large effect	Medium effect

Hypotheses

All hypotheses have been pre-registered through EGAP.

Hypothesis

*The ranking of the magnitudes of the decrease in harassment will be:
In-group/High status > $\frac{\text{In-group/Low status}}{\text{Out-group/High status}}$ > Out-group/Low status.*

Table: Experimental Design and Hypothesized Effect Sizes

	In-group	Out-group
Low followers	Medium effect	Small effect
High followers	Large effect	Medium effect

Hypotheses

All hypotheses have been pre-registered through EGAP.

Hypothesis

*The ranking of the magnitudes of the decrease in harassment will be:
In-group/High status > $\frac{\text{In-group/Low status}}{\text{Out-group/High status}}$ > Out-group/Low status.*

Table: Experimental Design and Hypothesized Effect Sizes

	In-group	Out-group
Low followers	Medium effect	Small effect
High followers	Large effect	Medium effect

Hypothesis

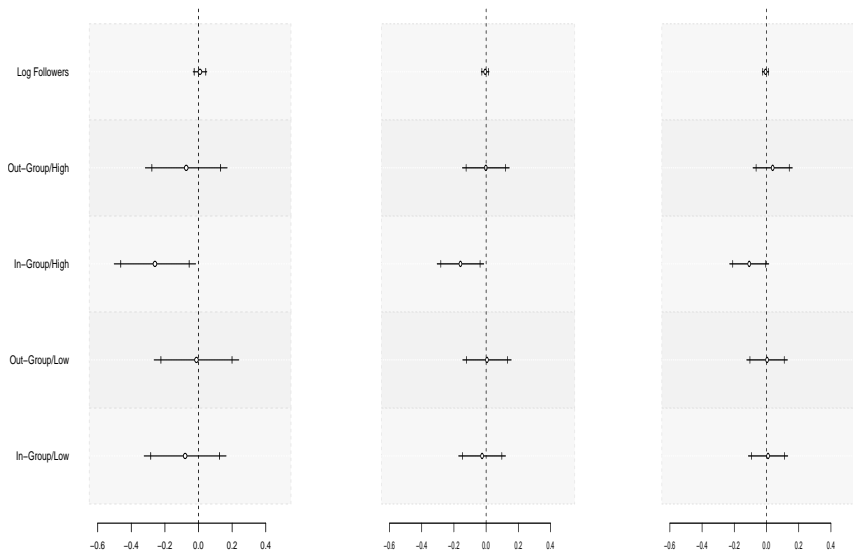
The decrease in offensive language will be smaller in subjects who provide less information in their profile.

Change in Racist Language: Full Sample (242)

1 Week

2 Weeks

1 Month



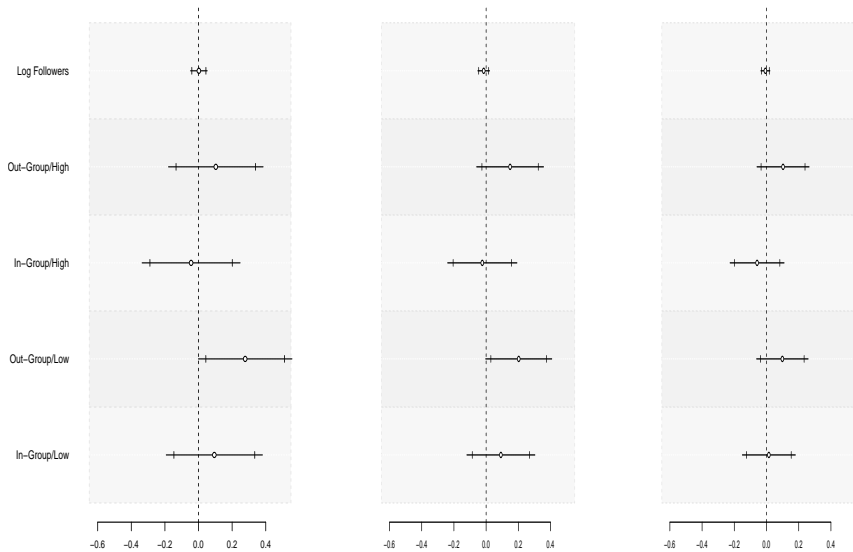
Change in Average Daily Slur Use

Change in Racist Language: Non-Anonymous Sample (84)

1 Week

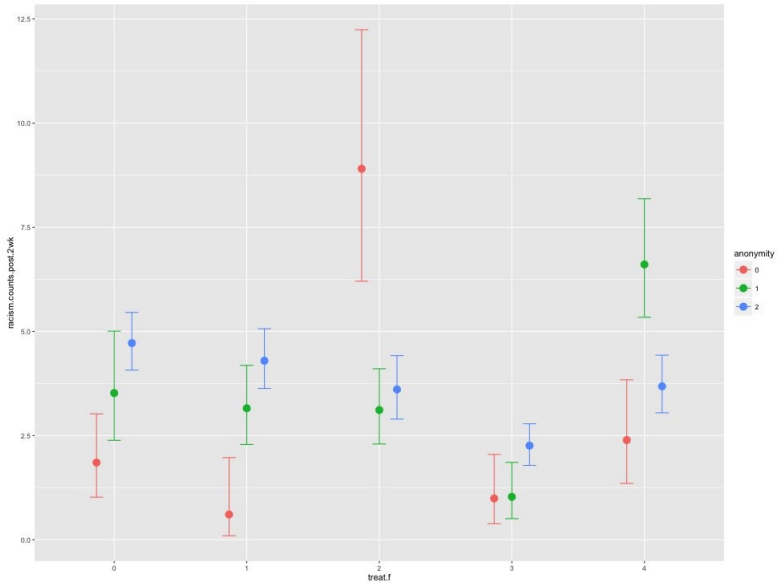
2 Weeks

1 Month



Change in Average Daily Slur Use

Science Moves Quickly



In Real World Terms

My intervention caused the 50 subjects in the most effective condition to tweet the word “nigger” an estimated 186 fewer times in the month after treatment.

Ongoing Project

- Political civility

Ongoing Project

- Political civility
- Incivility demobilizes and polarizes

A Visual Overview

 **Donald J. Trump** @realDonaldTrump · Oct 11
In Texas now, leaving soon for BIG rally in Florida

  8.4K  30K 

 **NSA agent jim** @ [redacted] · Oct 11
@realDonaldTrump HOPE U DIE


  6  92 


 **Sean** [redacted]  





@ [redacted] @realDonaldTrump typical libtard go fuck yourself

11:25 AM - 12 Oct 2016

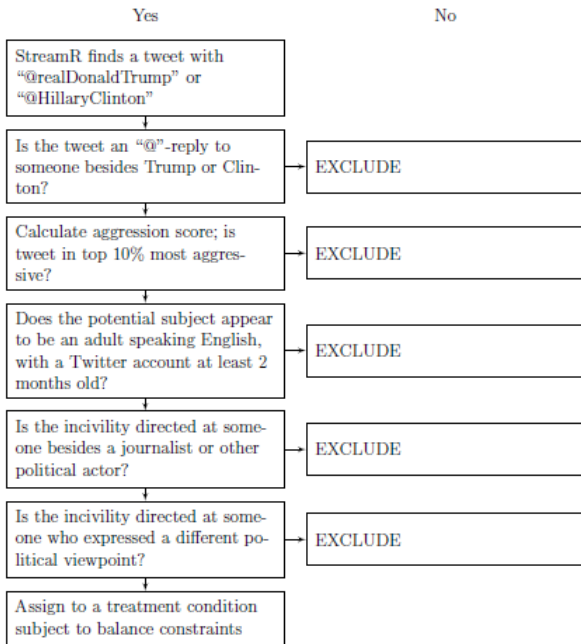
   

 [redacted] @realDonaldTrump

 **Matthew** [redacted] @Matthew [redacted] Oct 12 **FEELINGS**
@ [redacted] You shouldn't use language like that. Republicans need to remember that our opponents are real people, with real feelings.

Detection



Manipulations

- All bots are white men with many followers

Manipulations

- All bots are white men with many followers
- Trump-critical subjects get a tweet from a Democrat or Hillary bot

Manipulations

- All bots are white men with many followers
- Trump-critical subjects get a tweet from a Democrat or Hillary bot
- Test effectiveness of three types of messages

Manipulations

- All bots are white men with many followers
- Trump-critical subjects get a tweet from a Democrat or Hillary bot
- Test effectiveness of three types of messages
- Ideologically scale subjects (Barberá, 2015) and look for heterogeneous effects

Tweetment Variations

- Different rhetoric to appeal to different moral frameworks

Tweetment Variations

- Different rhetoric to appeal to different moral frameworks
 - ▶ Authority moral foundation: "You shouldn't use language like that. [Democrats/Republicans] need to behave according to the proper rules of political civility."

Tweetment Variations

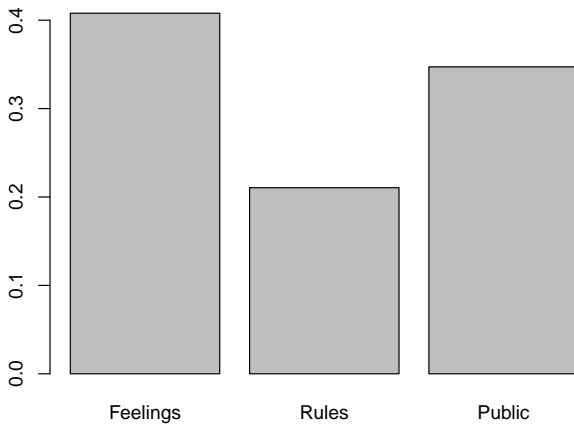
- Different rhetoric to appeal to different moral frameworks
 - ▶ Authority moral foundation: "You shouldn't use language like that. [Democrats/Republicans] need to behave according to the proper rules of political civility."
 - ▶ Care moral foundation: "You shouldn't use language like that. [Democrats/Republicans] need to remember that our opponents are real people, with real feelings."

Tweetment Variations

- Different rhetoric to appeal to different moral frameworks
 - ▶ Authority moral foundation: "You shouldn't use language like that. [Democrats/Republicans] need to behave according to the proper rules of political civility."
 - ▶ Care moral foundation: "You shouldn't use language like that. [Democrats/Republicans] need to remember that our opponents are real people, with real feelings."
 - ▶ Placebo message: "Remember that everything you post here is public. Everyone can see that you tweeted this."

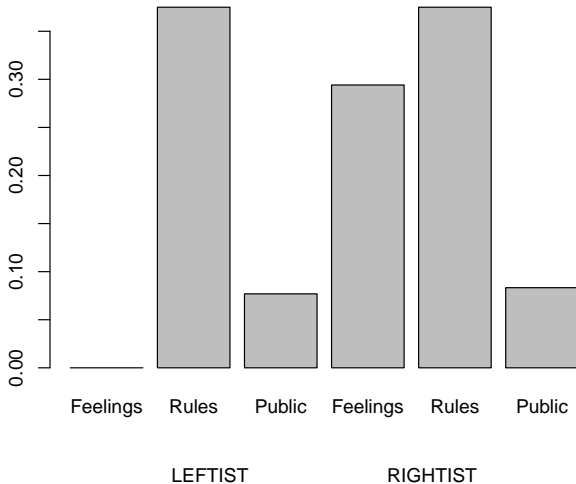
Preliminary Results

Response Rates by Treatment (N=224)



Preliminary Results

Percentage of Conciliatory Response (N=72)



An Optimistic Note



Veteran [redacted] @ [redacted] · 11h

@realDonaldTrump a nasty woman beat the piss out of you in front of millions, yet u claim victory. You are a bad hombre



62



437



1st Sgt. [redacted] @ [redacted] · 11h

@ [redacted] @realDonaldTrump repeat that statement when ure head gets sawed off.



3



Veteran [redacted] @ [redacted] · 11h

@ [redacted] @realDonaldTrump u know what's absolutely hilarious, is I've patrolled Mosul and dealt w al queda snipers. I don't fear shit



37



An Optimistic Note



1st Sgt. [redacted] · 11h

@ [redacted] @realDonaldTrump I spilled blood in the desert for two tours one in storm and one in shield, ure point is.



1st Sgt. [redacted] · 11h

@ [redacted] @realDonaldTrump its ure.vote. ure right. I just don't see the logic. Good luck



Veteran [redacted] · 11h

@ [redacted] @realDonaldTrump I'm happy to argue, but not to be called a disgrace. I'm not on day zero of basic



1st Sgt. [redacted] · 11h

@ [redacted] @realDonaldTrump ure opinion is ures and mine is mine. Your right tho I should of kept the fact of ure disgrace to myself



Veteran [redacted] · 11h

@ [redacted] @realDonaldTrump i don't really take it personally if it's another vet. I'm sorry i was upset



An Optimistic Note

 1st Sgt. [redacted] · 11h
@ [redacted] @realDonaldTrump sorry for the comment
← ↻ ❤️ 2 ⋮

 Veteran [redacted] · 11h
@ [redacted] @realDonaldTrump me too. I have mad respect for you, one tour was enough
← ↻ ❤️ 4 ⋮

 1st Sgt. [redacted] · 11h
@ [redacted] @realDonaldTrump My grandfather , same division. He had it tougher though. Omaha Beach ww2
← ↻ ❤️ 1 ⋮

 Veteran [redacted] · 11h
@ [redacted] @realDonaldTrump jesus, yes.
← ↻ ❤️ ⋮

 1st Sgt. [redacted] · 11h
@ [redacted] @realDonaldTrump These wars are gonna be the longest and most difficult to win. Its like a widespread Vietnam
← ↻ ❤️ 2 ⋮

 Veteran [redacted] · 11h
@ [redacted] @realDonaldTrump absolutely. 15 years and counting, way too much. At least there is some hope in Iraq
← ↻ ❤️ 2 ⋮



1st Sgt. [redacted]



Follow

@ [redacted] @realDonaldTrump With the right leader and Generals in place

LIKE

1



Thanks for your comments, and for listening!

- km2713@nyu.edu
- @kmmunger (no harassment, please)

Attrition rates

	Control	A	B	C	D
Baseline # of subjects	40	49	44	50	48
# with > 1 Post-treatment tweets	40	46	42	47	47
# with > 25 Post-treatment tweets	40	34	33	35	43
Attrition %, < 25 tweets	10%	18%	16%	18%	4%

The number of subjects who tweeted more than 1 or 25 times after the application of the treatment.