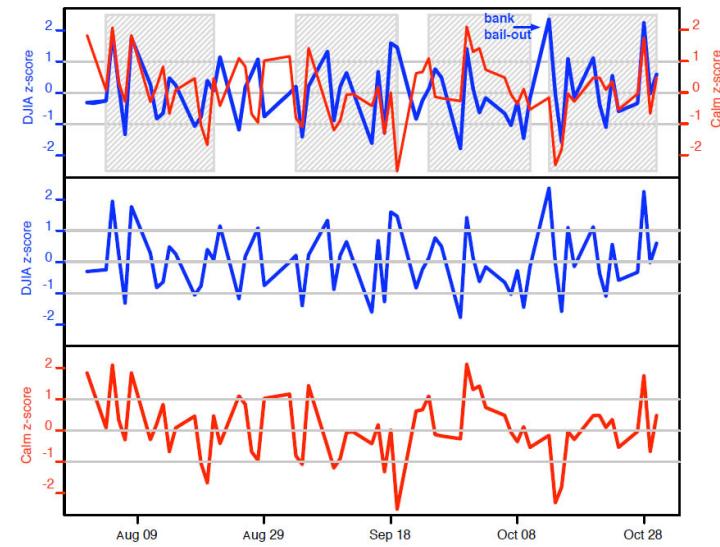


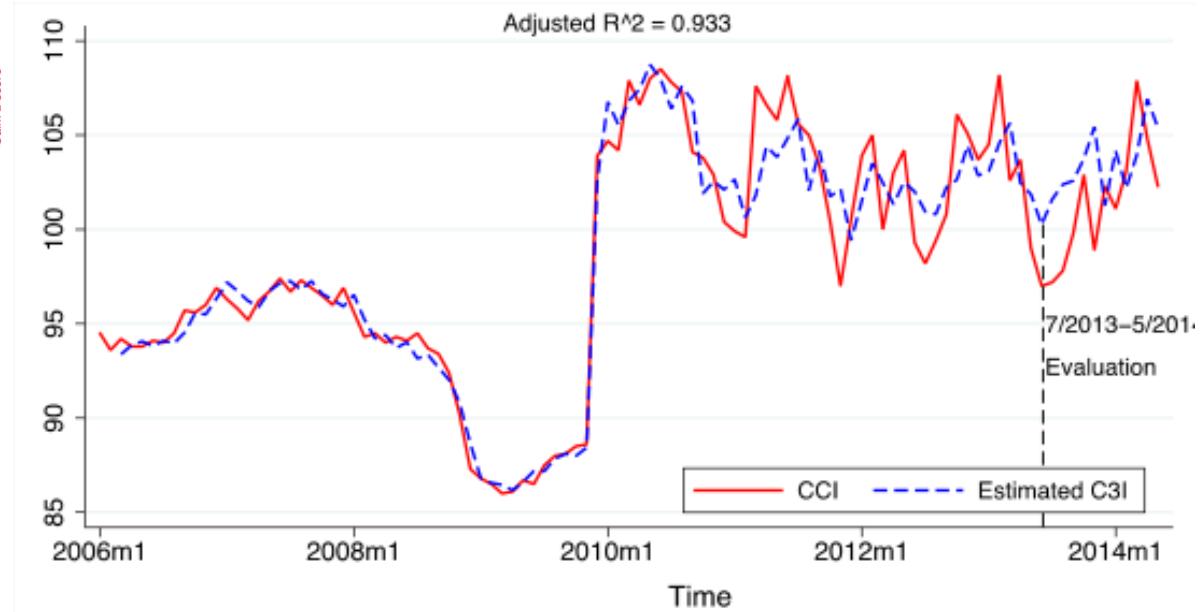
Using social media indicators to study regional socio-economic resilience

Johan Bollen & Krishna Bathina

Modeling and predicting socio-economic phenomena from large-scale online data



Bollen et al (2011) Predicting
market returns from social media
sentiment



Dong & Bollen (2014) Modeling Chinese
consumer confidence from online search data

Catch a black swan

Statistical regularities of the relation between online indicators and socio-economic phenomena.

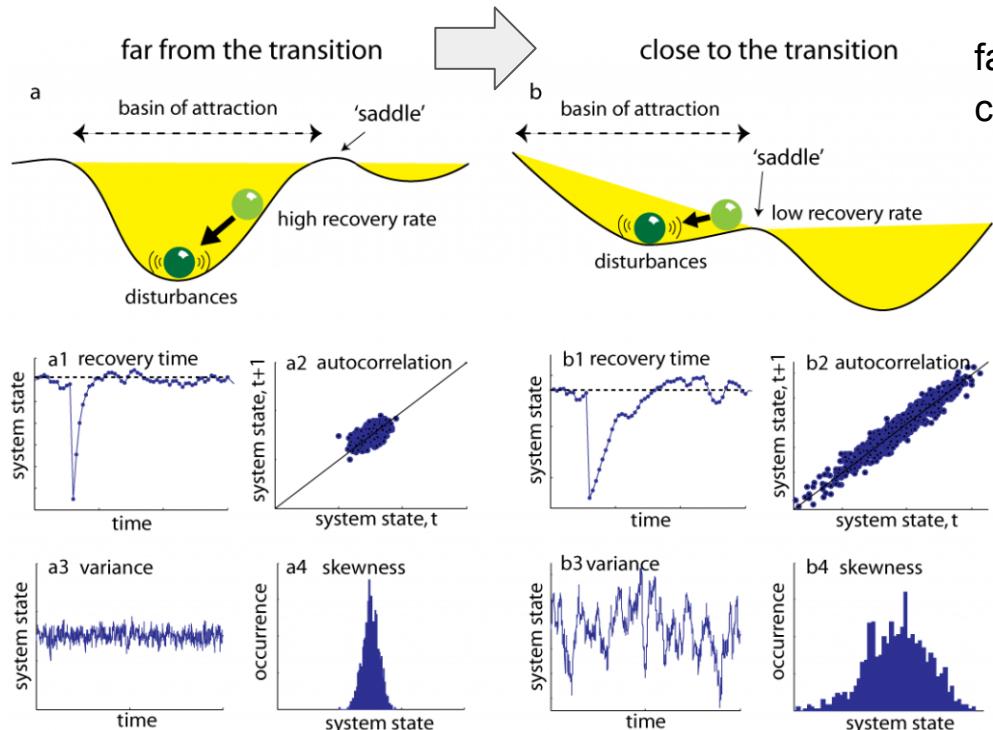
BUT, many socio-technical systems undergo rapid, difficult to predict, and infrequent transitions:

- Economic collapse
- Community collapse
- Social order collapse

=> may be more important to predict than overall dynamics



Early warning indicators of critical transitions



far from transition = High resilience = fast recovery
close to transition = low resilience = slow recovery

- Early warning indicators:
- Increasing auto-correlation
 - Variance

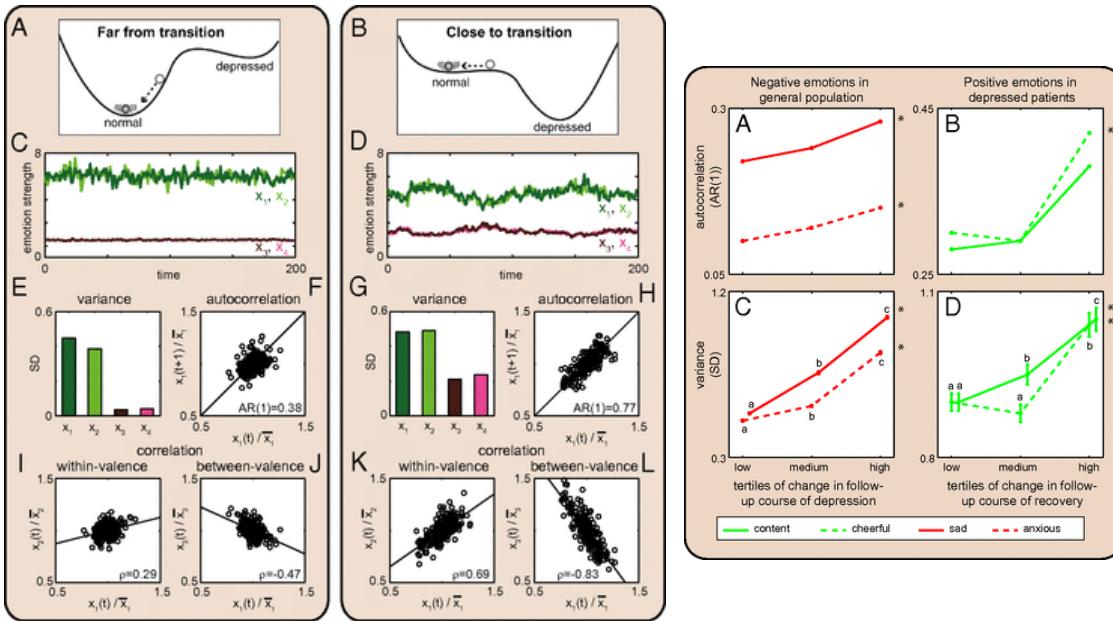
Quantitative early warning indicators from longitudinal dynamics of system parameters

<http://www.early-warning-signals.org/>

Scheffer et al (2012) Anticipating Critical Transitions. Science, 338(6105):344-348

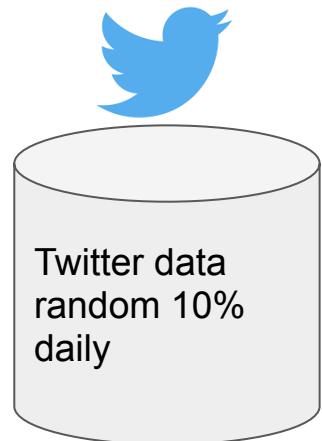
Our approach

1. Large-scale social media data
2. Temporal patterns that express social and economic resilience for a region
3. Early warning indicators of potential transition

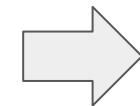


Van de Leemput et al (2014) **Critical slowing down as early warning for the onset and termination of depression**. PNAS January 7, 2014 111 (1) 87-92

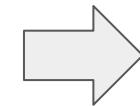
Data and Processing



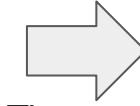
Indiana University
Network institute



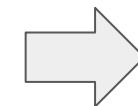
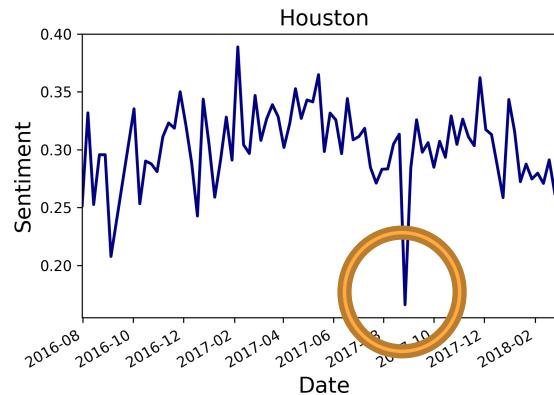
Lon/Lat



Content
analysis



Time
stamp



Regional
early warning
indicators

"favorited": false,

CONTENT

"text": "One of my favorite Tony Gillian\u0027s!! #SexyVillian #MichealDrucker #MCM #MCE #tonygoldwyn @... <https://t.co/l3lxtyMRYZ>",

"created_at": "Mon Aug 01 17:15:45 -0400 2016",

"user_id_str": "77403942",

TIME STAMP

"user_screen_name": "Gladiator6082",

"retweet_count": 0,

"retweeted": false,

"source": "\u003ca href\u003d\"http://instagram.com\"\nrel\u003d\"nofollow\"\u003eInstagram\u003c/a\u003e",

"id_str": "760222511644680193",

**"entities": {
 "hashtags": [...],
 "urls": [...],
},**

**"coordinates": {
 "type": "Point",
 "coordinates": [
 -85.8445,
 31.3275
]
},**

GEO-LOCATION

"user": {...},

"place": {...}

VADER - Valence Aware Dictionary and sEntiment Reasoner

- Best performer in large-scale benchmark (Ribeiro et al, 2016)
- Crowd-sourced valence lexicon (7,516 terms)
- 5 grammatical and syntactical rules geared towards social media

1. **Punctuation** (!) modifies intensity: “great” vs. “great!!!!”
2. **Capitalization** (all caps): “great” vs “**GREAT**”
3. **Degree Modifiers** (adjectives/adverbs): “That was **really** great”
4. **Contrastive Conjunction** (but): “That was fun **but** I didn’t like it”
5. **Trigram** analysis to find negation: “That was not that great”

“**Johan is smart, handsome, and funny.**”

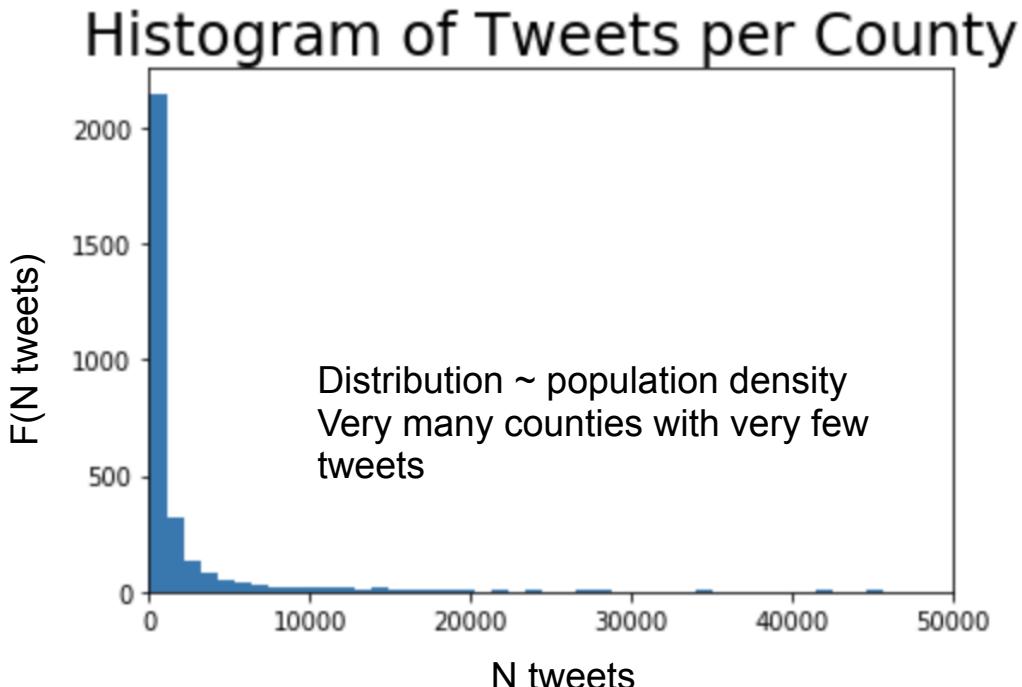
- Neg: 0.0, Neu: 0.254, Pos: 0.746
- **Compound: +0.8316**

“**Today SUX!**”

- Neg: 0.779, Neu: 0.221, Pos: 0.0
- **Compound: -0.5461**

Our data

- 3,221 US counties (US Census Bureau 2010)
- 11,063,495 total tweets
- Mean 3434.8 Tweets/county
- 19 month period

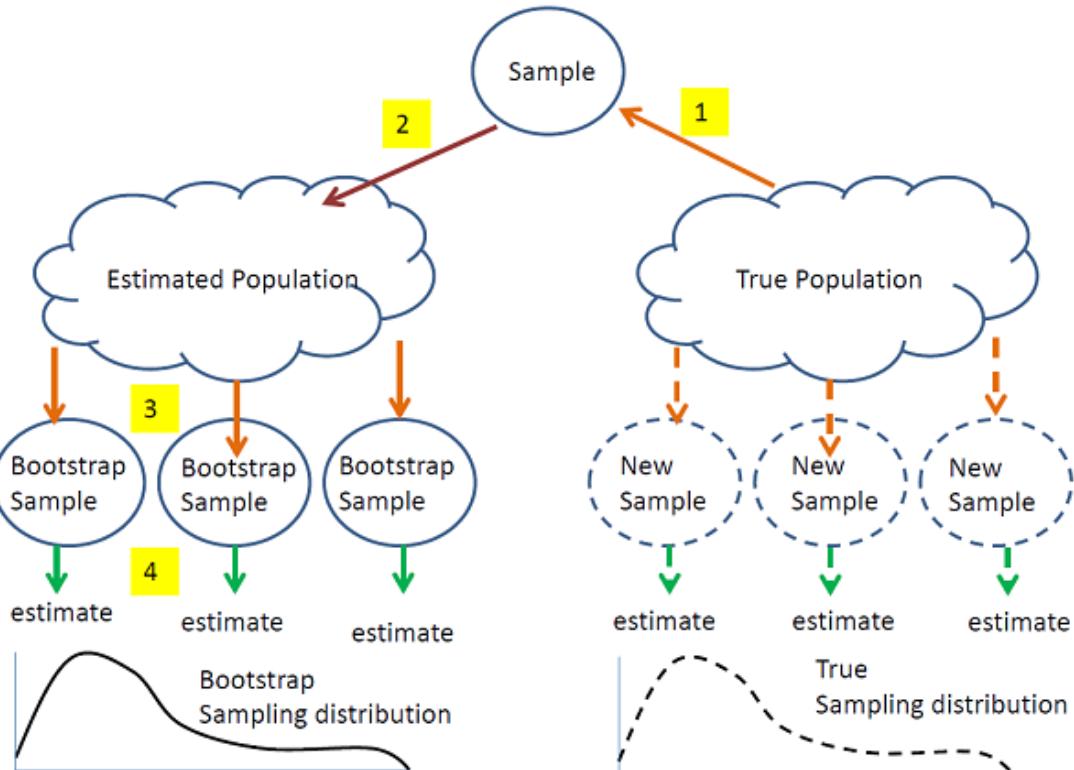


Most counties will have few tweets, increasing uncertainty of our estimate of “average sentiment” at time t

Bootstrapping time series

Rationale: Twitter posts at irregular intervals, so different samples for each time window.
Uncertainty estimated mean/time?

- Bootstrap Twitter sample for each time period by randomly sampling with replacement
- 95% CIs expresses uncertainty



Face validity

Testing our regional sentiment signal for three cities

- **Florida (Sep 2017)**
- **Houston (Aug 2017)**
- **Puerto Rico (Sep 2017)**

Significant sentiment change in areas when afflicted by hurricane
AND, counter-factual: no sentiment change in unaffected areas?

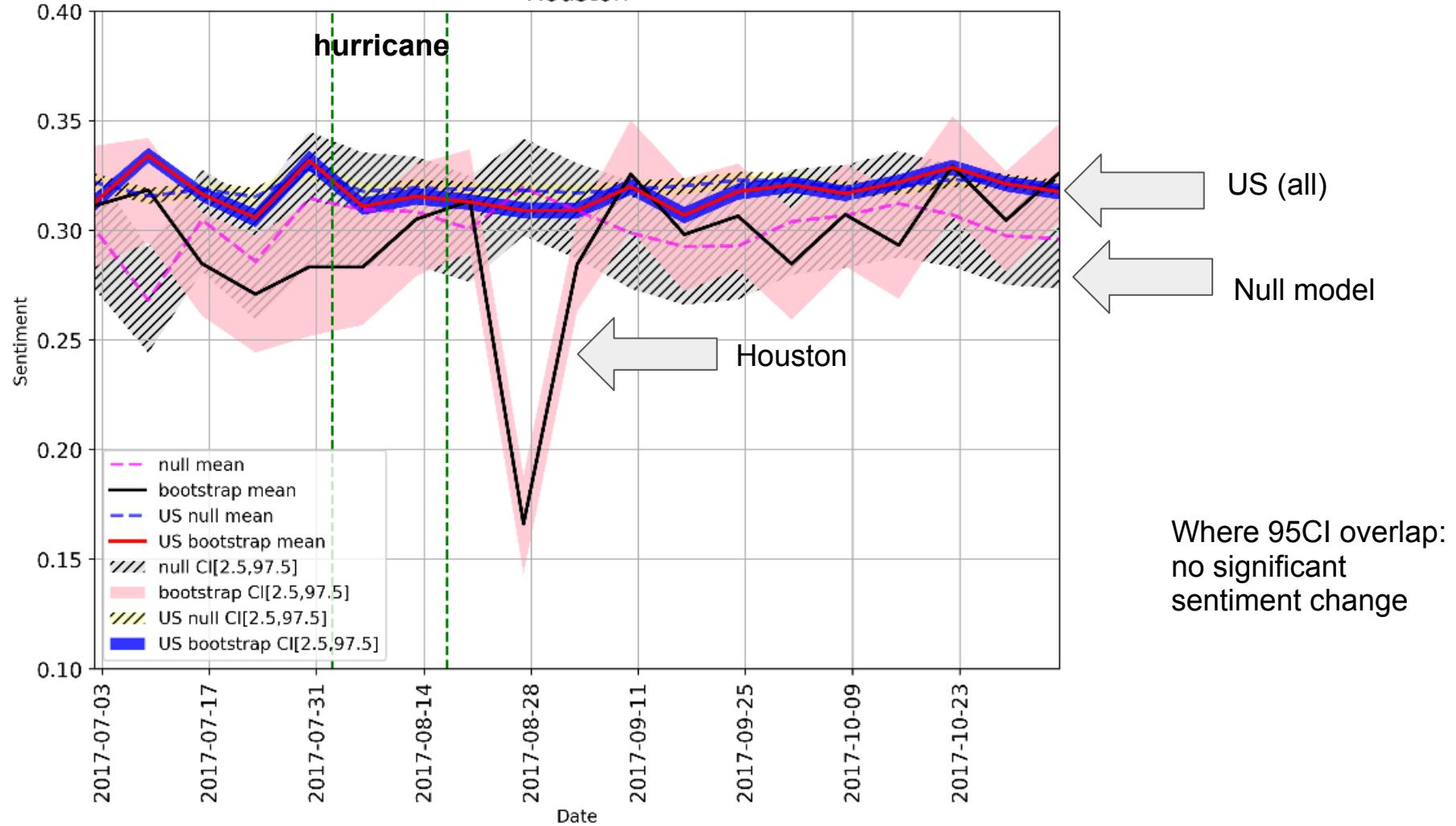
Compare sentiment levels at time t to appropriate null-models for same time t:

- (1) bootstrap sentiment time series to estimate 95% CI (N=10,000)
- (2) US baseline (all counties)
- (3) Null-model: random selection of tweets with same "day of week distribution"

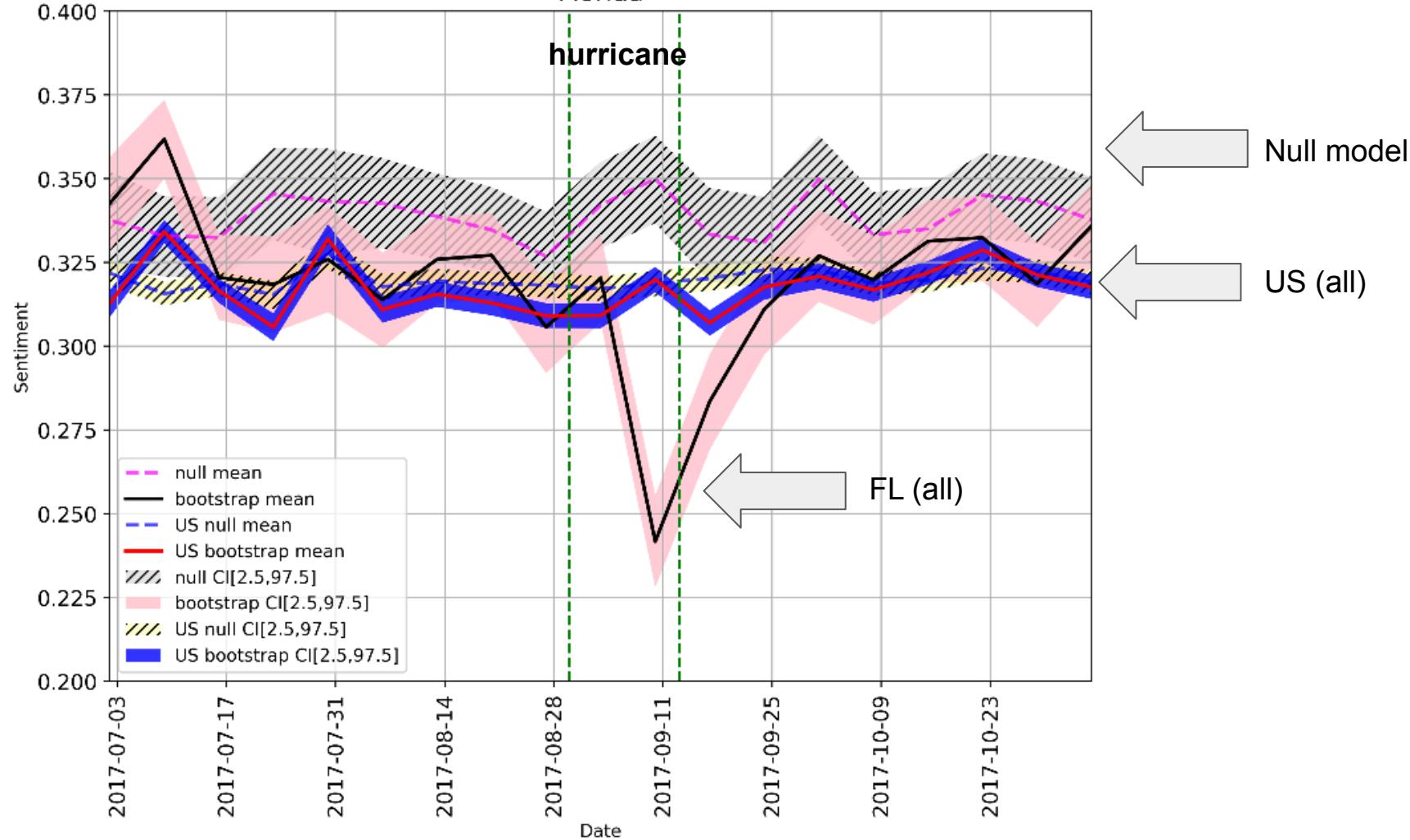
Su	M	T	W	Th	F	S
20	22	30	30	35	45	59

Null-model follows same distribution to mitigate weekly cycles

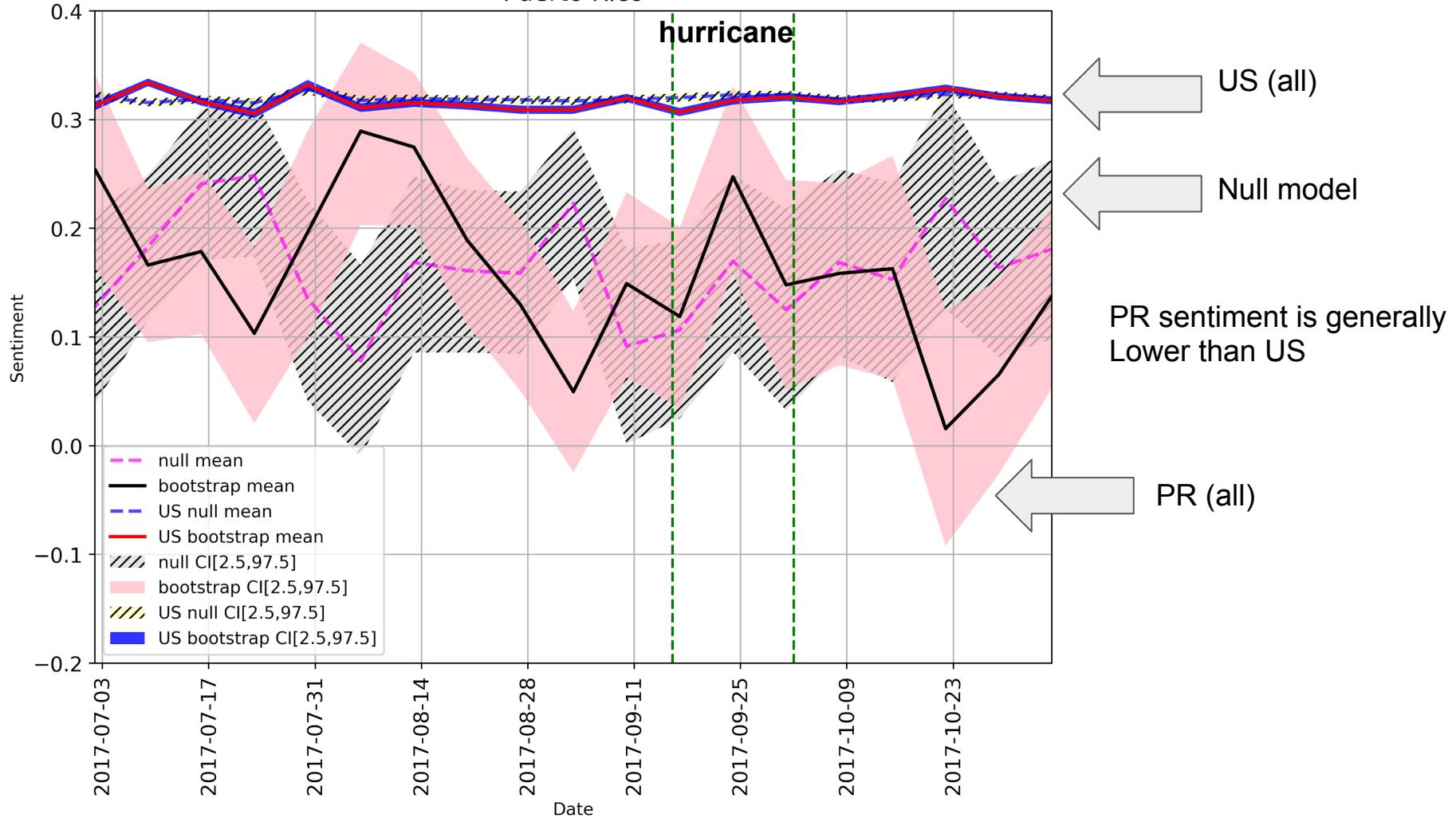
Houston

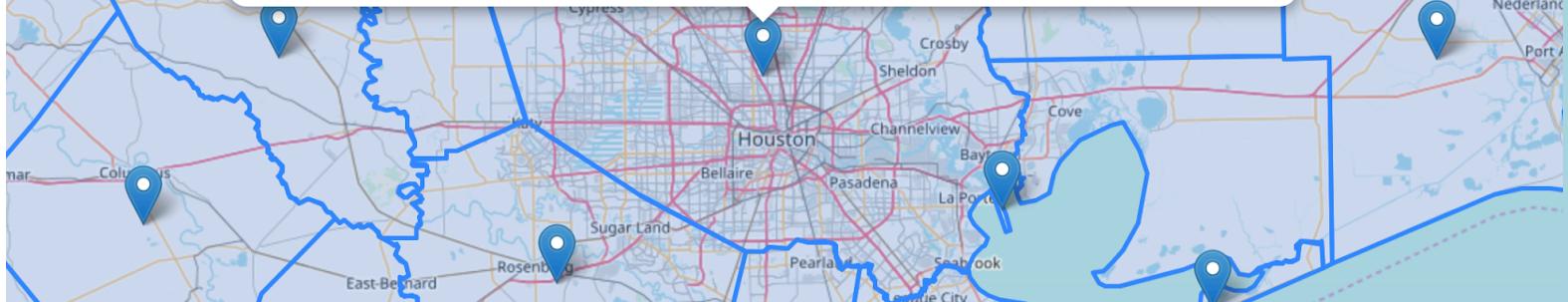
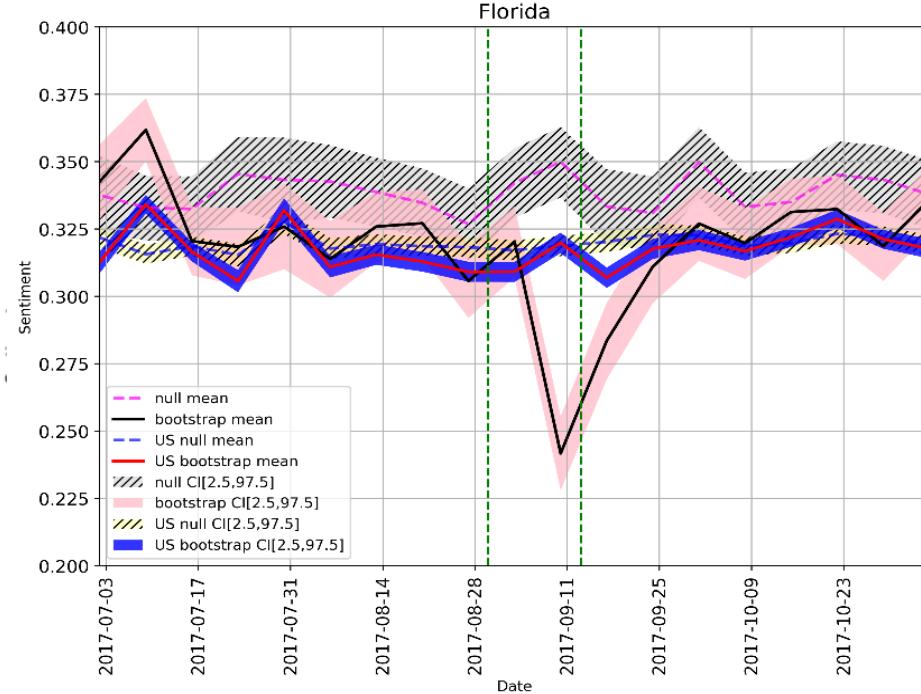


Florida



Puerto Rico





Beyond valence

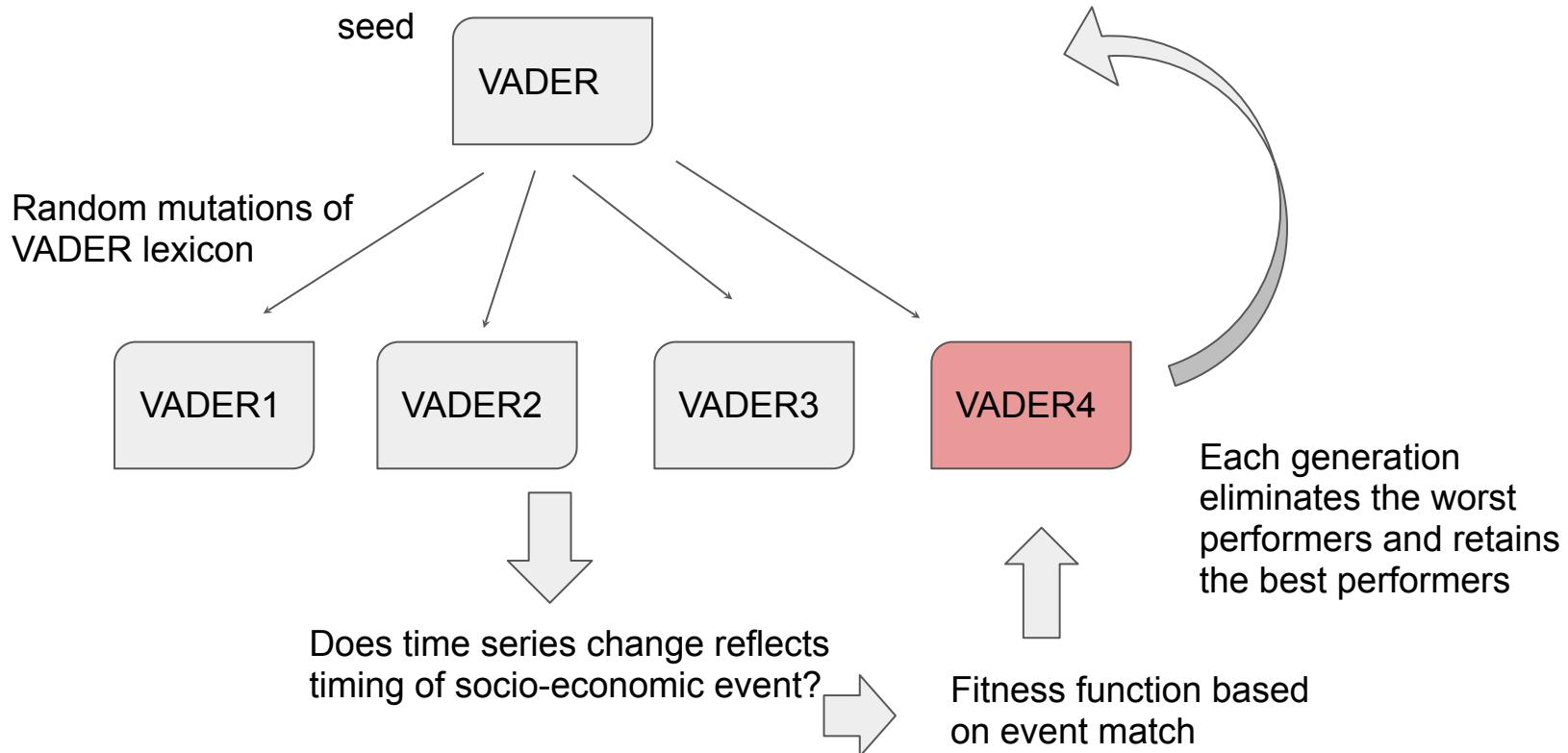
Valence is very broad category of sentiment, but how about...

- Culture and collective personality
 - Specific socio-economic trends, e.g. unrest, dissatisfaction, apprehension, confidence

General approach:

- Pick specific event with **ground truth of event occurrence at given time t**, e.g Ferguson unrest
 - “Breed” a NLP tool whose lexicon is specifically suited to generate the most specific signal with respect to the occurrence of that socio-economic phenomenon in time



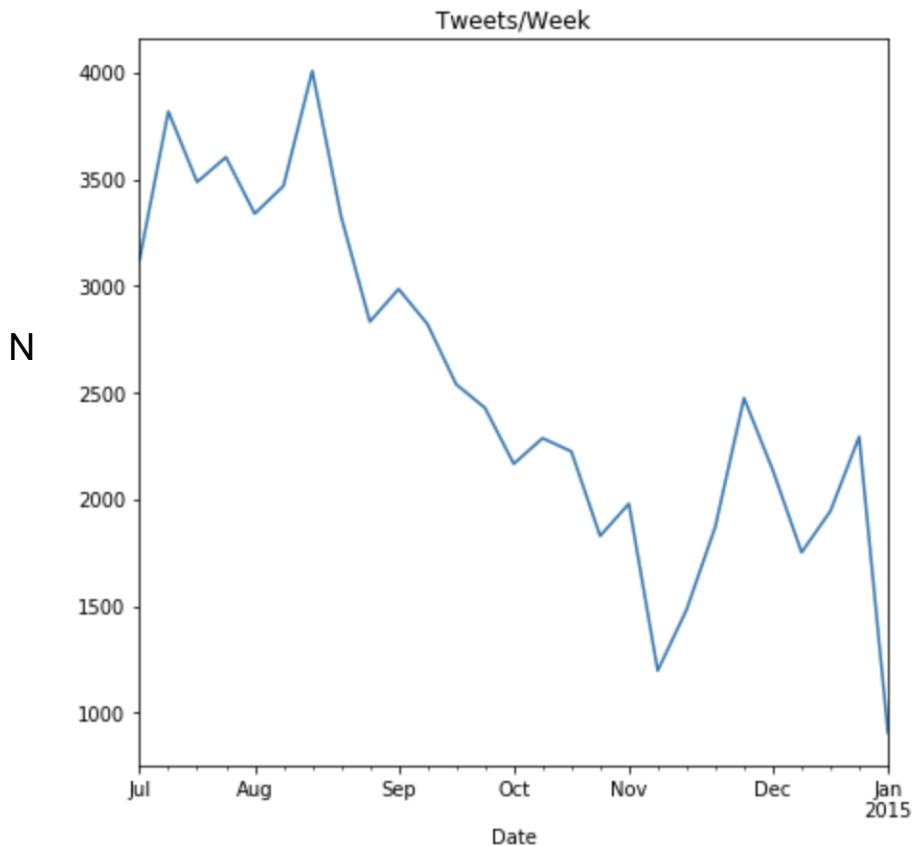


Breeding a “social unrest” indicator

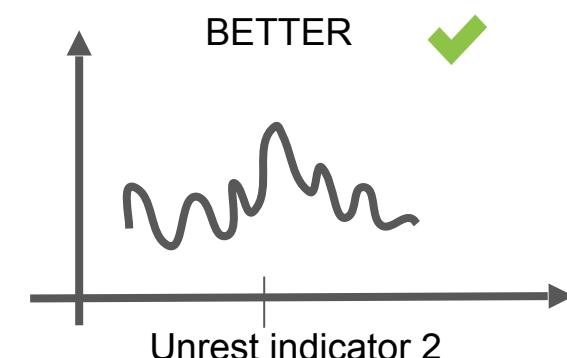
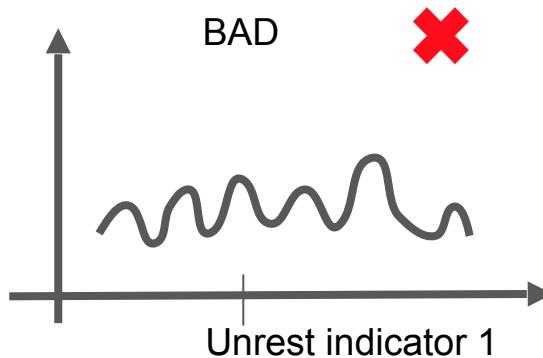
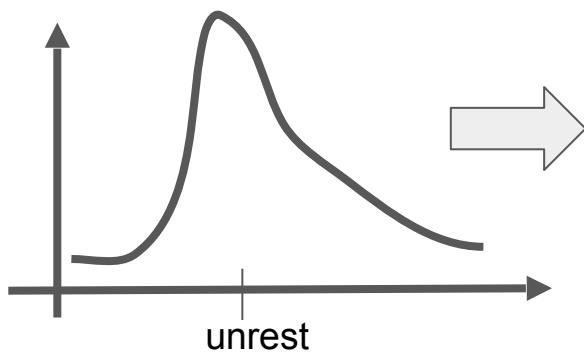
Ferguson MO riots in 2016

- All tweets in St. Louis county from July - December 2016
- N = 68,333 Tweets
- Mean weekly resamples
- Sentiment converted to Z-scores

$$z = \frac{x - \mu}{\sigma}$$



Fitness: does our indicator respond to social unrest?

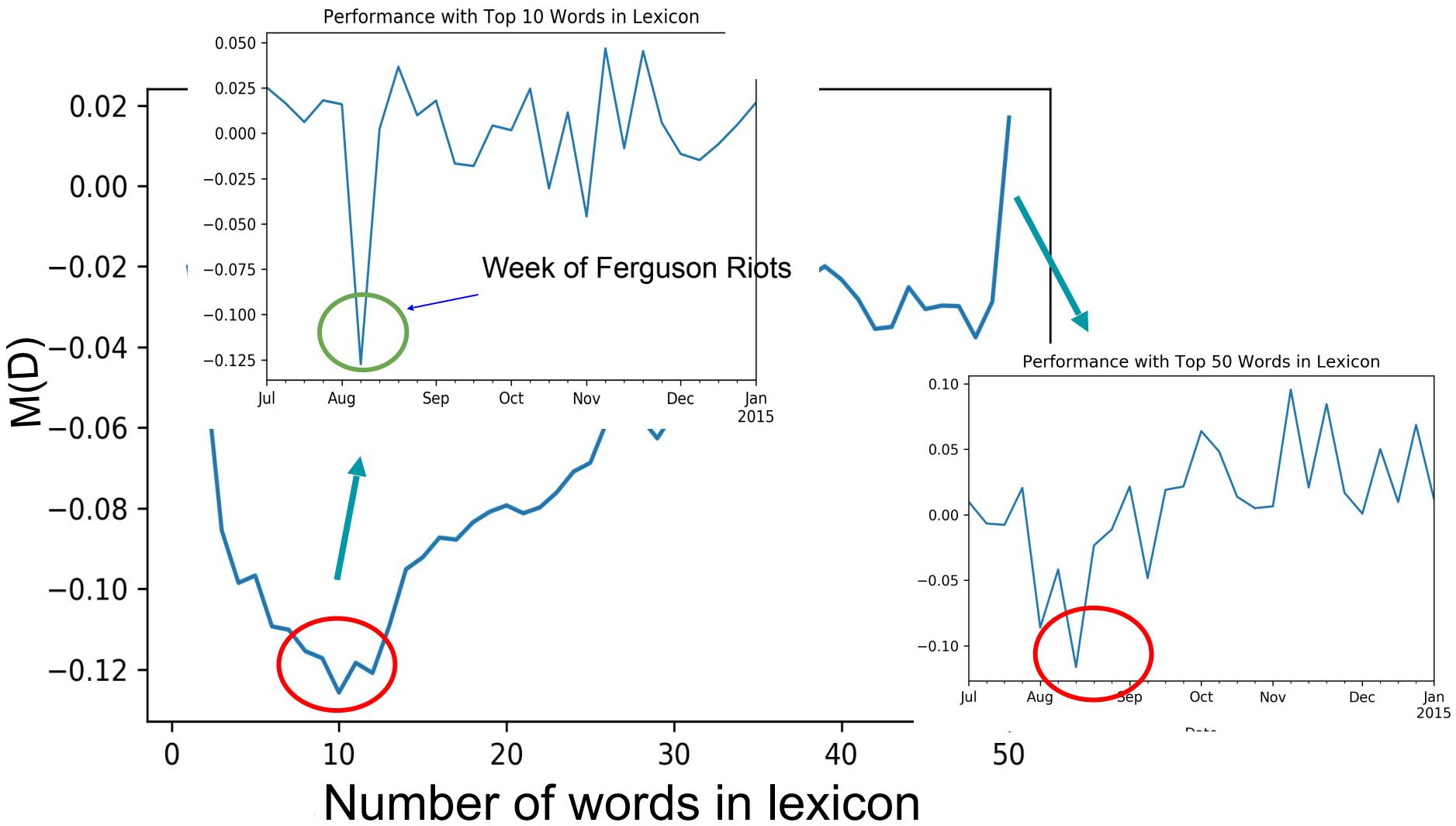


$$M(D) = \mathcal{H} + \min(D.z) + \frac{|\arg \min_{D.w} D.z - D.w_R|}{|D.w_F - D.w_R|}$$

Binned Entropy
(specificity)

+ Max value
(response magnitude)

+ Ratio of distance from unrest
of minimum value compared
to the farthest difference



Continuing work

We have the time series! Now:

- Train sentiment analyzer for recognition of other “black swans”: riots, emergencies, disasters, economic cycles
- Detection of critical transitions: early warning indicators?
- Contribution to EDA indicator of regional social and economic resilience

Acknowledgement: Marten Scheffer and Ingrid van de Leemput, Wageningen U Economic Development Agency, Timothy Slaper of Indiana U