Geographical variation of happiness as expressed by users of Twitter

Kameron Decker Harris

Peter S. Dodds Isabel M. Klouman Catherine Bliss Christopher M. Danforth





April 26, 2011

1 Methodology

- Valence scoring
- Geotagging

2 Results

- Maps
- Wordshifts
- Socioeconomic correlations

3 Conclusions

Why measure happiness?

It's important:

- Choose: gross domestic... product or happiness?
- People ≠ rational ... feelings influence everything. To model behavior, we must understand states of mind.

Why measure happiness?

It's important:

- Choose: gross domestic... product or happiness?
- People ≠ rational ... feelings influence everything. To model behavior, we must understand states of mind.

Why measure happiness?

It's important:

- Choose: gross domestic... product or happiness?
- People ≠ rational ... feelings influence everything. To model behavior, we must understand states of mind.

Why measure happiness?

It's important:

- Choose: gross domestic... product or happiness?
- People ≠ rational ... feelings influence everything. To model behavior, we must understand states of mind.

Why measure happiness?

It's important:

- Choose: gross domestic... product or happiness?
- People ≠ rational ... feelings influence everything. To model behavior, we must understand states of mind.

Hedonom-*a-what*?



(actually GDP!)

'hedonometer' = device to measure well-being

Ideally:

- objective
- improvable
- text-based
- fast, for big data
 (Twitter data set has ~ 30 × 10⁹ words)

Hedonom-*a-what*?



(actually GDP!)

- objective
- improvable
- text-based
- fast, for big data (Twitter data set has \sim 30 imes 10⁹ words

Hedonom-*a-what*?



(actually GDP!)

- objective
- improvable
- text-based
- fast, for big data
 (Twitter data set has \sim 30 × 10⁹ words

Hedonom-*a-what*?



(actually GDP!)

- objective
- improvable
- text-based
- fast, for big data
 (Twitter data set has \sim 30 × 10⁹ words)

Hedonom-*a-what*?



(actually GDP!)

'hedonometer' = device to measure well-being Ideally:

- objective
- improvable
- text-based
- fast, for big data

(Twitter data set has \sim 30 imes 10^9 words)

Hedonom-*a-what*?



(actually GDP!)

- objective
- improvable
- text-based
- \blacksquare fast, for big data (Twitter data set has $\sim 30 \times 10^9$ words)

-Valence scoring

ANEW = 'Affective Norms for English Words'



1034 words with scored on three dimensions:

- valence (happiness)
- arousal
- dominance

from Bradley and Lang (1999)

-Valence scoring

ANEW = 'Affective Norms for English Words'



- 1034 words with scored on three dimensions:
 - valence (happiness)
 - arousal
 - dominance

from Bradley and Lang (1999)

-Valence scoring

ANEW = 'Affective Norms for English Words'



- 1034 words with scored on three dimensions:
 - valence (happiness)
 - arousal
 - dominance
- from Bradley and Lang (1999)

- Methodology
 - Valence scoring

Scoring a text

- look from occurances of words in ANEW list, construct frequencies f_i, i = 1,...,1034
- create weighted average based on scores *s_i*

happiness
$$=rac{\sum_i s_i f_i}{\sum_i f_i}$$

Valence scoring

Scores for typical ANEW words



triumphant (8.82) / paradise (8.72) / love (8.72) luxury (7.88) / trophy (7.78) / glory (7.55) optimism (6.95) / church (6.28) / pancakes (6.08) street (5.22) / paper (5.20) / engine (5.20) neurotic (4.45) / vanity (4.30) / derelict (4.28) fault (3.43) / lawsuit (3.37) / corrupt (3.32) disgusted (2.45) / hostage (2.20) / trauma (2.10) funeral (1.39) / rape (1.25) / suicide (1.25)

Takeaway: ANEW sample a mixture of salient and neutral words

- Methodology
 - └─Valence scoring

Example scores for certain texts

Text:	h_{avg}	Words with a similar score:
Soul/Gospel music	6.9	chocolate (6.88), leisurely
lyrics [20]		(6.88), penthouse (6.81)
Pop music lyrics [20]	6.7	dream (6.73) , honey (6.73) ,
		sugar (6.74)
Dante's Paradise [33]	6.5	muffin (6.57) , rabbit (6.57) ,
		smooth (6.58)
Tweets, $9/9/2008$ to	6.4	thought (6.39) , face (6.39) ,
12/31/2010 (present		blond (6.42)
work)		
Rock music lyrics [20]	6.3	church (6.28) , tree (6.32) , air
		(6.34)
Enron Emails [34]	6.2	clouds (6.18), alert (6.20),
		computer (6.24)
State of the Union	6.1	grass (6.12) , idol (6.12) , bottle
Messages [20]		(6.15)
New York Times	6.0	hotel (6.00) , tennis (6.02) ,
(1987–2007) [35]		wonder (6.03)
Blogs [20]	5.8	owl (5.80), whistle (5.81),
		humble (5.86)
Dante's Inferno [33]	5.5	glacier (5.50) , repentant
		(5.53), mischief (5.57)
Metal/Industrial	5.4	lamp (5.41) , elevator (5.44) ,
music lyrics [20]		truck (5.47)

Kameron Decker Harris

-Valence scoring

So many possibilities!

One can examine:

- trends in time
- network characteristics
- beyond ANEW (distributions of all words, information theory)
- spatial trends

Valence scoring

So many possibilities!

One can examine:

trends in time

- network characteristics
- beyond ANEW (distributions of all words, information theory)

spatial trends

Valence scoring

So many possibilities!

One can examine:

- trends in time
- network characteristics
- beyond ANEW (distributions of all words, information theory)

spatial trends

Valence scoring

So many possibilities!

One can examine:

- trends in time
- network characteristics
- beyond ANEW (distributions of all words, information theory)

spatial trends

Valence scoring

So many possibilities!

One can examine:

- trends in time
- network characteristics
- beyond ANEW (distributions of all words, information theory)
- spatial trends

Valence scoring

So many possibilities!

One can examine:

- trends in time
- network characteristics
- beyond ANEW (distributions of all words, information theory)
- **spatial trends** \leftarrow this talk

- Methodology
 - Geotagging

Two routes:

1 self-reporting

- users accounts specify a "Location", e.g. "Burlington VT", "New York", "over the rainbow!! :-)"
- parse out city & state, e.g. ["burlington", "vt"]
- look up in database of places (USGS Geonames)
- \sim ~5-20% identifiable this way
- store at level of county
- data since Sept. 2008
- we do this now

- Methodology
 - Geotagging

Two routes:

self-reporting

- users accounts specify a "Location", e.g. "Burlington VT", "New York", "over the rainbow!! :-)"
- parse out city & state, e.g. ["burlington", "vt"]
- look up in database of places (USGS Geonames)
- \sim 5-20% identifiable this way
- store at level of county
- data since Sept. 2008
- we do this now

- Methodology
 - Geotagging

- Two routes:
 - self-reporting
 - users accounts specify a "Location", e.g. "Burlington VT", "New York", "over the rainbow!! :-)"
 - parse out city & state, e.g. ["burlington", "vt"]
 - look up in database of places (USGS Geonames)
 - \sim ~5-20% identifiable this way
 - store at level of county
 - data since Sept. 2008
 - we do this now
 - 2 location encoded in tweet

- Methodology
 - Geotagging

- Two routes:
 - self-reporting
 - users accounts specify a "Location", e.g. "Burlington VT", "New York", "over the rainbow!! :-)"
 - parse out city & state, e.g. ["burlington", "vt"]
 - look up in database of places (USGS Geonames)
 - \sim ~5-20% identifiable this way
 - store at level of county
 - data since Sept. 2008
 - we do this now
 - 2 location encoded in tweet

- Methodology
 - Geotagging

- Two routes:
 - self-reporting
 - users accounts specify a "Location", e.g. "Burlington VT", "New York", "over the rainbow!! :-)"
 - parse out city & state, e.g. ["burlington", "vt"]
 - look up in database of places (USGS Geonames)
 - \sim -5-20% identifiable this way
 - store at level of county
 - data since Sept. 2008
 - we do this now
 - 2 location encoded in tweet

- Methodology
 - Geotagging

- Two routes:
 - self-reporting
 - users accounts specify a "Location", e.g. "Burlington VT", "New York", "over the rainbow!! :-)"
 - parse out city & state, e.g. ["burlington", "vt"]
 - look up in database of places (USGS Geonames)
 - $\blacksquare~\sim\!\!5\text{--}20\%$ identifiable this way
 - store at level of county
 - data since Sept. 2008
 - we do this now
 - 2 location encoded in tweet

- Methodology
 - Geotagging

- Two routes:
 - self-reporting
 - users accounts specify a "Location", e.g. "Burlington VT", "New York", "over the rainbow!! :-)"
 - parse out city & state, e.g. ["burlington", "vt"]
 - look up in database of places (USGS Geonames)
 - \blacksquare $\sim\!\!5\text{--}20\%$ identifiable this way
 - store at level of county
 - data since Sept. 2008
 - we do this now

- Methodology
 - Geotagging

- Two routes:
 - self-reporting
 - users accounts specify a "Location", e.g. "Burlington VT", "New York", "over the rainbow!! :-)"
 - parse out city & state, e.g. ["burlington", "vt"]
 - look up in database of places (USGS Geonames)
 - \blacksquare $\sim\!\!5\text{--}20\%$ identifiable this way
 - store at level of county
 - data since Sept. 2008
 - we do this now

- Methodology
 - Geotagging

- Two routes:
 - self-reporting
 - users accounts specify a "Location", e.g. "Burlington VT", "New York", "over the rainbow!! :-)"
 - parse out city & state, e.g. ["burlington", "vt"]
 - look up in database of places (USGS Geonames)
 - \blacksquare $\sim\!\!5\text{--}20\%$ identifiable this way
 - store at level of county
 - data since Sept. 2008
 - we do this now

- Methodology
 - Geotagging

- Two routes:
 - self-reporting
 - users accounts specify a "Location", e.g. "Burlington VT", "New York", "over the rainbow!! :-)"
 - parse out city & state, e.g. ["burlington", "vt"]
 - look up in database of places (USGS Geonames)
 - \blacksquare $\sim\!\!5\text{--}20\%$ identifiable this way
 - store at level of county
 - data since Sept. 2008
 - we do this now
 - 2 location encoded in tweet

-									-
(-eog	ranhical	variation	ot	hanniness	as ev	pressed	hv	lisers of	witter
OCOF	rapincar	variation		nappiness,		pressea	ωγ.	03013 01	WILLEL

0 1		
— Results		
L _{Maps}		

The big picture: maps

L_Maps

Counties with data for > 100 unique users ('09-'10)

Happiness in US by County



All states ('09-'10)

Happiness in US by State



Geographical variation of happiness, as expressed by users of Twitter

Results

└─ Wordshifts

Getting into the details: word-shift plots

└_ Wordshifts

Why states' scores differ: specific words



VT with respect to the U.S.



└_ Wordshifts

Why states' scores differ: specific words



CA with respect to the U.S.



└_ Wordshifts

Why states' scores differ: specific words



DC with respect to the U.S.



Geographical variation of happiness, as expressed by users of Twitter

Results

-Socioeconomic correlations

Spatial correlations with socioeconomic indices

Socioeconomic correlations

Happiness and \$\$\$



Socioeconomic correlations

Happiness, information, and politics



Conclusions, future directions

Fully investigate socioeconomic correlations

- Add in time analysis (deal with non-uniform sample sizes)
- Compare to detailed geocoding
- Kridgings, other elegant geostatistical analysis

Conclusions, future directions

- Fully investigate socioeconomic correlations
- Add in time analysis (deal with non-uniform sample sizes)
- Compare to detailed geocoding
- Kridgings, other elegant geostatistical analysis

Conclusions, future directions

- Fully investigate socioeconomic correlations
- Add in time analysis (deal with non-uniform sample sizes)
- Compare to detailed geocoding
- Kridgings, other elegant geostatistical analysis

Conclusions, future directions

- Fully investigate socioeconomic correlations
- Add in time analysis (deal with non-uniform sample sizes)
- Compare to detailed geocoding
- Kridgings, other elegant geostatistical analysis

Extras

Happiness trajectories of individual states



Extras

Information trajectories of individual states

