# Social Media Text Analysis

Stony Brook University
CSE545, Fall 2016

# Basics of Natural Language Processing

- Tokenization
  - Sentence
  - Word

- Part of Speech Tagging

- Syntactic Parsing

# From language to features

Feature encodings

- Count
- Relative Frequency
- TF-IDF

- Dimensionally Reduced

# Features: Closed-to-Open Vocabulary

# Standard Tasks

- Insight


- Prediction

# General "Insight" Framework

# Prediction Framework

# Levels of Analysis

# Example Tasks

1. Text-based Geolocation

2. Community Health Prediction
   (Handling many features, few observations)

3. Human Temporal Orientation
   (Sophisticated Features)

# 1. **Text-based Geolocation**

GOAL: Determine where a given user lives.

Versions

1. Based on posts (e.g. status updates, tweets)
2. Based on profile information

Gold-Standard: Geo-coordinates (lat+lon)

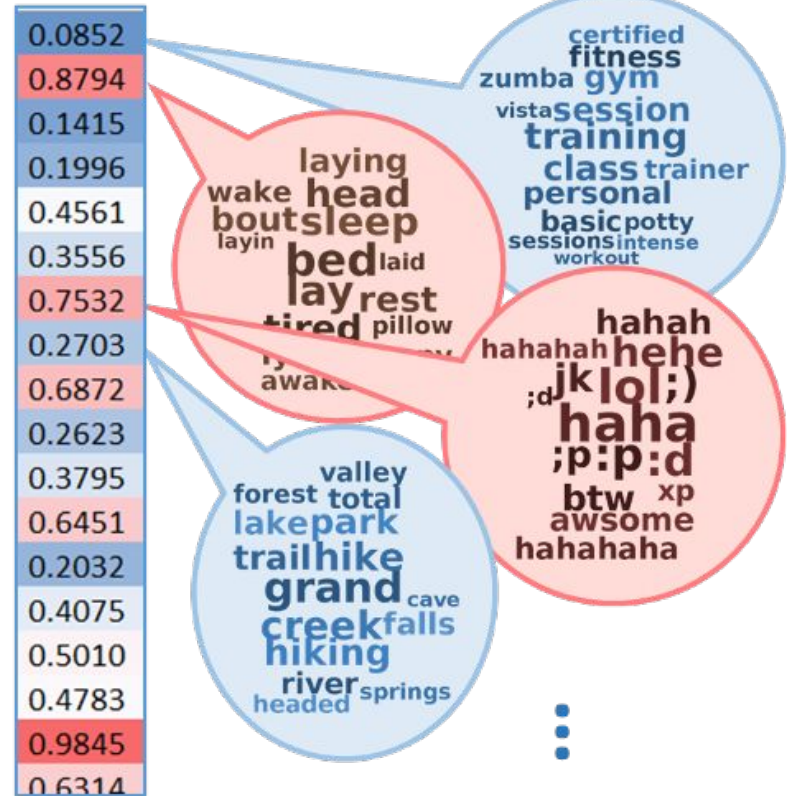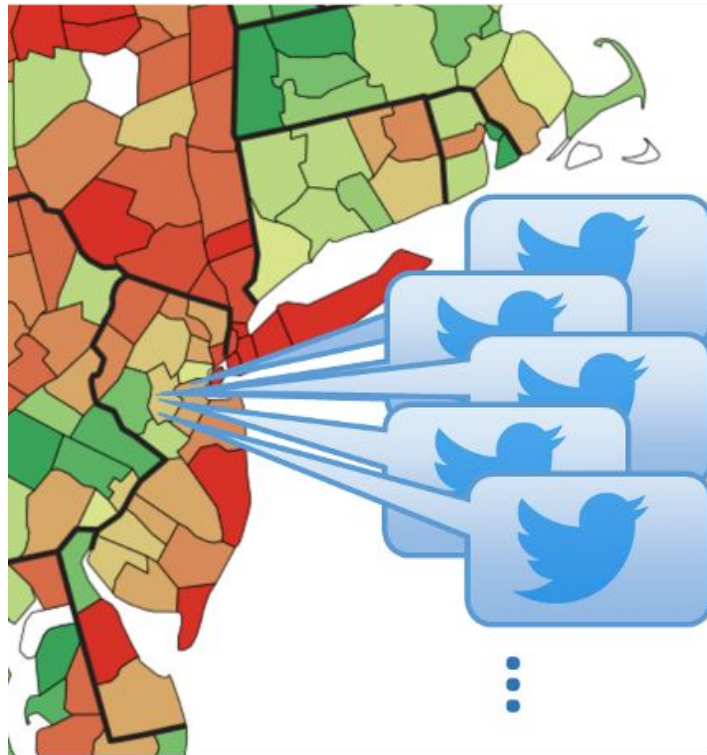# 2. Community Health Prediction

## Data



Atherosclerotic heart disease mortality

# Encoding a community

# Twitter Predicts Heart Disease

Eichstaedt, J. C., Schwartz, H. A., Kern, M. L., Park, G.,..., Ungar, L. H., & Seligman, M. E. (2015). Psychological Language on Twitter Predicts County-Level Heart Disease Mortality. *Psychological Science* 26(2), 159-169

# 3. Human Temporal Orientation

# Building a model

| message | R1 | R2 | R3 | m | class |
|---|---|---|---|---|---|
| *did nothing this morning but watch TV and it was fantastic =)* | -.67 | -.50 | -.50 | -.55 | past |
| *dislikes being sick.... and misses her bf* | 0 | 0 | 0 | 0 | present |
| *pancake day tomorrow pancake day tomorrow xxxxx* | .50 | .50 | 1 | .67 | future |

Training Data
4.3k
tweets+ statuses

→ Learn Model → Model

Application Data
1.3m statuses

# Building a model

| message | R1 | R2 | R3 | m | class |
|---|---|---|---|---|---|
| *did nothing this morning but watch TV and it was fantastic =)* | -.67 | -.50 | -.50 | -.55 | past |
| *dislikes being sick.... and misses her bf* | 0 | 0 | 0 | 0 | present |
| *pancake day tomorrow pancake day tomorrow xxxxx* | .50 | .50 | 1 | .67 | future |

Linguistic Feature Extraction

# Building a model

| message | | R1 | R2 | R3 | m | class |
|---------|---|----|----|----|----|-------|
| *did p...* | *...it was fantastic =)* | | | | | *past* |
| *dis...* | | | | | | *...sent* |
| *pancake day tomorrow pancake day tomorrow xxxxx* | | .50 | .50 | 1 | .67 | future |

parts-of-speech
(covers tense)

time
expressions

Linguistic Feature Extraction

lexica

words and
phrases

# Building a model

| message | | R1 | R2 | R3 | m | class |
|---|---|---|---|---|---|---|
| did n... | ...d it was fantastic =) | | | | | past |
| dis... | | | | | | ...sent |
| pancake day tomorrow pancake day tomorrow xxxxx | | .50 | .50 | 1 | .67 | future |

"today"

"in two weeks"

parts-of-speech
(covers tense)

time
expressions

"January 15"

"last year"

Linguistic Feature Extraction

lexica

words and
phrases

# Building a model

parts-of-speech (covers tense)

time expressions

Linguistic Feature Extraction

lexica

words and phrases

# Building a model

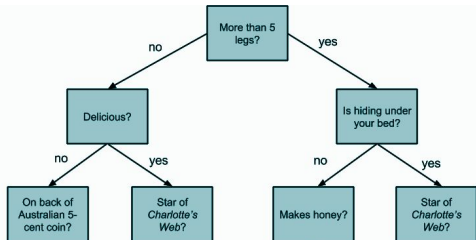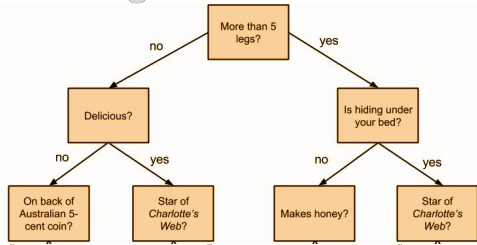| message | R1 | R2 | R3 | m | class |
|---|---|---|---|---|---|
| *did nothing this morning but watch TV and it was fantastic =)* | -.67 | -.50 | -.50 | -.55 | past |
| *dislikes being sick.... and misses her bf* | 0 | 0 | 0 | 0 | present |
| *pancake day tomorrow pancake day tomorrow xxxxx* | .50 | .50 | 1 | .67 | future |

Linguistic Feature Extraction

Learn Message-Level Model

# Building a model



## Accuracy over a held-out set: 72%;   baseline: 53%

Schwartz, H. A., Park, G., Sap, M., ..., & Ungar, L. (2015). Extracting Human Temporal Orientation from Facebook Language. *NAACL-2015: Conference of the North American Chapter of the Association for Computational Linguistics*

# Building a model

| message | R1 | R2 | R3 | m | class |
|---|---|---|---|---|---|
| *did nothing this morning but watch TV and it was fantastic =)* | -.67 | -.50 | -.50 | -.55 | past |
| *dislikes being sick.... and misses her bf* | 0 | 0 | 0 | 0 | present |
| *pan........row xxxxx* | | | | | future |

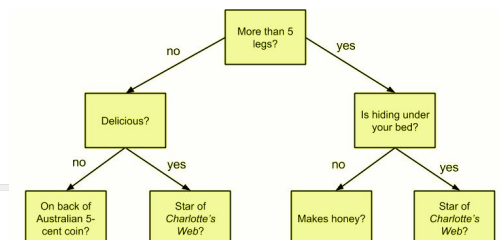parts-of-speech (covers tense) 62%

time expressions 59%

Linguistic Feature Extraction

lexica 68%

words and phrases 69%

...am Message-Level...

Accuracy over a held-out set: 72%;    baseline: 53%

Schwartz, H. A., Park, G., Sap, M., ..., & Ungar, L. (2015). Extracting Human Temporal Orientation from Facebook Language. *NAACL-2015: Conference of the North American Chapter of the Association for Computational Linguistics*

Legend: past (green), present (orange), future (blue)

Vertical axis: $r$ with gridlines at 0.2, 0.1, 0, -0.1, -0.2

Categories: age, gender, conscientiousness, impulsiveness, openness, extraversion, agreeableness, neuroticism

Apply to Participant Messages