# Linguistic Temporal Trajectory Analysis on Video Transcripts

Maximilian Mozes

Workshop on Linguistic Temporal Trajectory Analysis

2018 European Symposium Series on Societal Challenges in Computational Social Science

December 05, 2018

Kleinberg, B., Mozes, M. and van der Vegt, I., 2018. **Identifying the sentiment styles of YouTube's vloggers.**

In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing* (pp. 3581 - 3590).

# Video blogs

- Blogs in video format
- People filming their (daily) activities
- Can be domain-specific, e.g.
  - Technical product reviews
  - Beauty vlogs
  - „How-to" vlogs

# Continuous sentiment

- Vloggers try to arouse viewer's interest
- Sentiment as a means to achieve that?
- We measure continuous sentiment in videos
- Clustering approach to group similar sentiment styles

# Why video transcripts?

- Large amounts of data
- „Implicit annotations"
- Use of language in videos

# Why vlogs?

Many **samples** from single source

↓

Many **sources** for the same domain

↓

Many **domains**

Daily life vlogs, technical reviews, beauty vlogs, ...

# Data

- **27,333** vlog transcripts (24 users (13 male, 3 female))
- Each transcript consists of „textual chunks", e.g.

| | |
|---|---|
| 1 | *there are so many boogers in my nose* |
| 2 | *right now* |
| 3 | *forgot my memory card of my blog camera* |
| 4 | *in my room so now we're starting the* |
| 5 | *vlog on my phone what's going on I am so* |
| 6 | *not awake right now my makeup is* |
| 7 | *actually a hot and disaster* |

**Problem: no punctuations!**

# Dealing with non-punctuated data

- Sentiment analysis and non-punctuated data do not work well together

- Solution: analyze sentiment's neighborhood and check for
  - **negators** („not", „never")
  - **(de-)amplifiers** („really", „hardly")
  - **adversative conjunctions** („but", „yet")

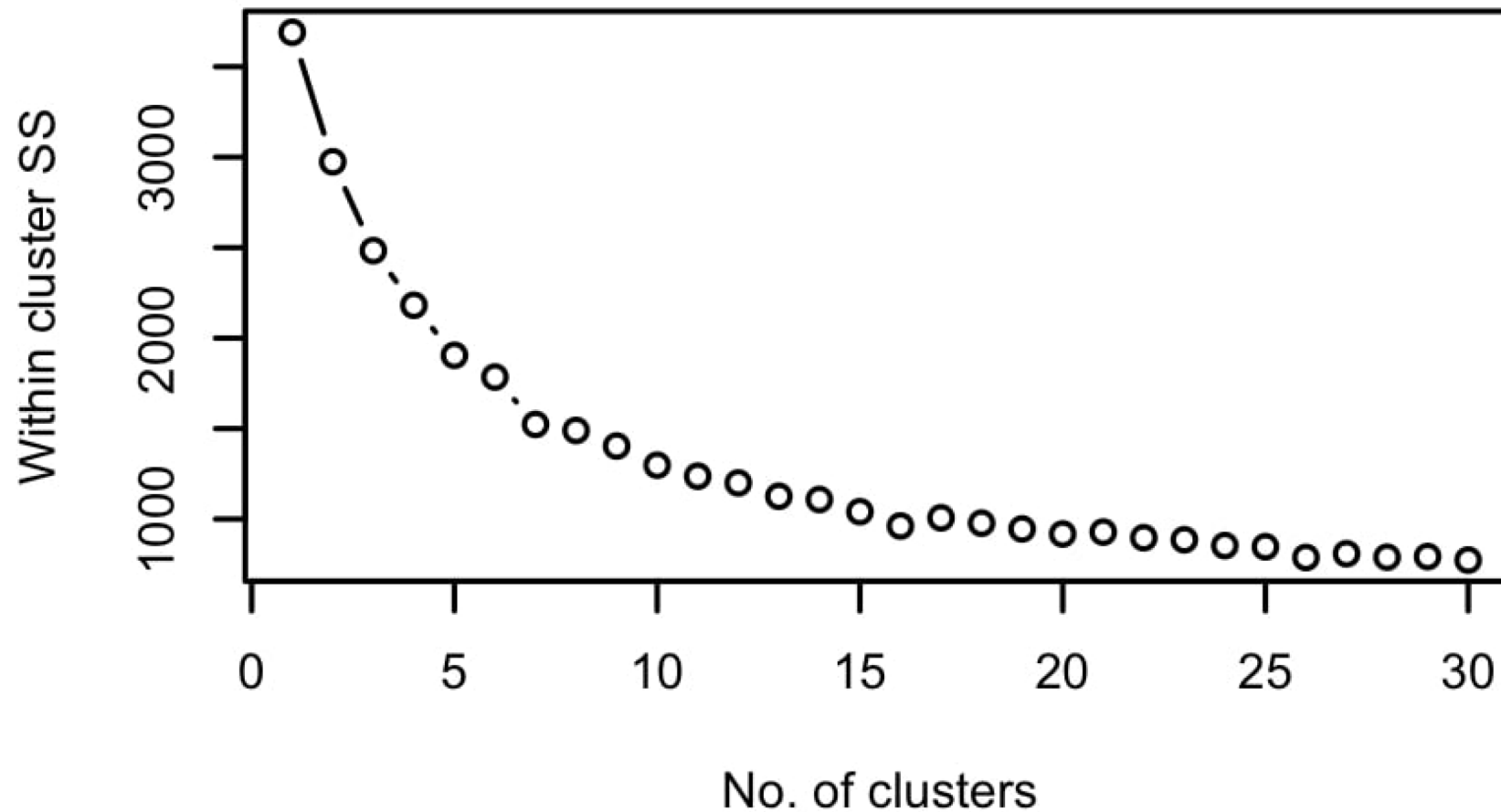„...this was not a **bad** day at all ..."

-3          +3

# Pipeline

1. **Identify sentiment values** for each transcript (Jockers & Rinker Polarity Lookup Table (Rinker, 2018))

2. **Normalize** sentiment values to 100-dimensional vector

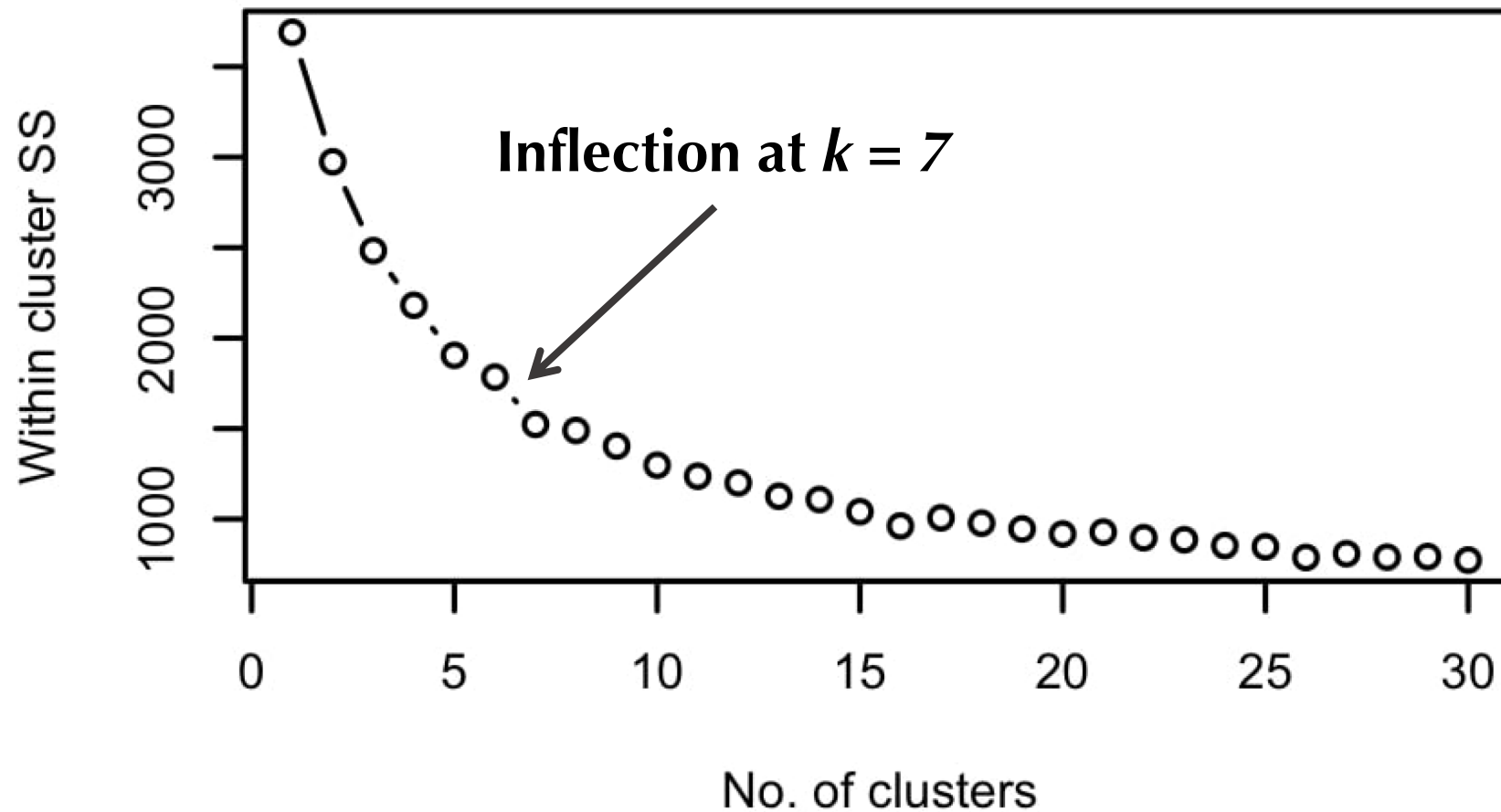3. *k*-means clustering to identify groups of sentiment styles
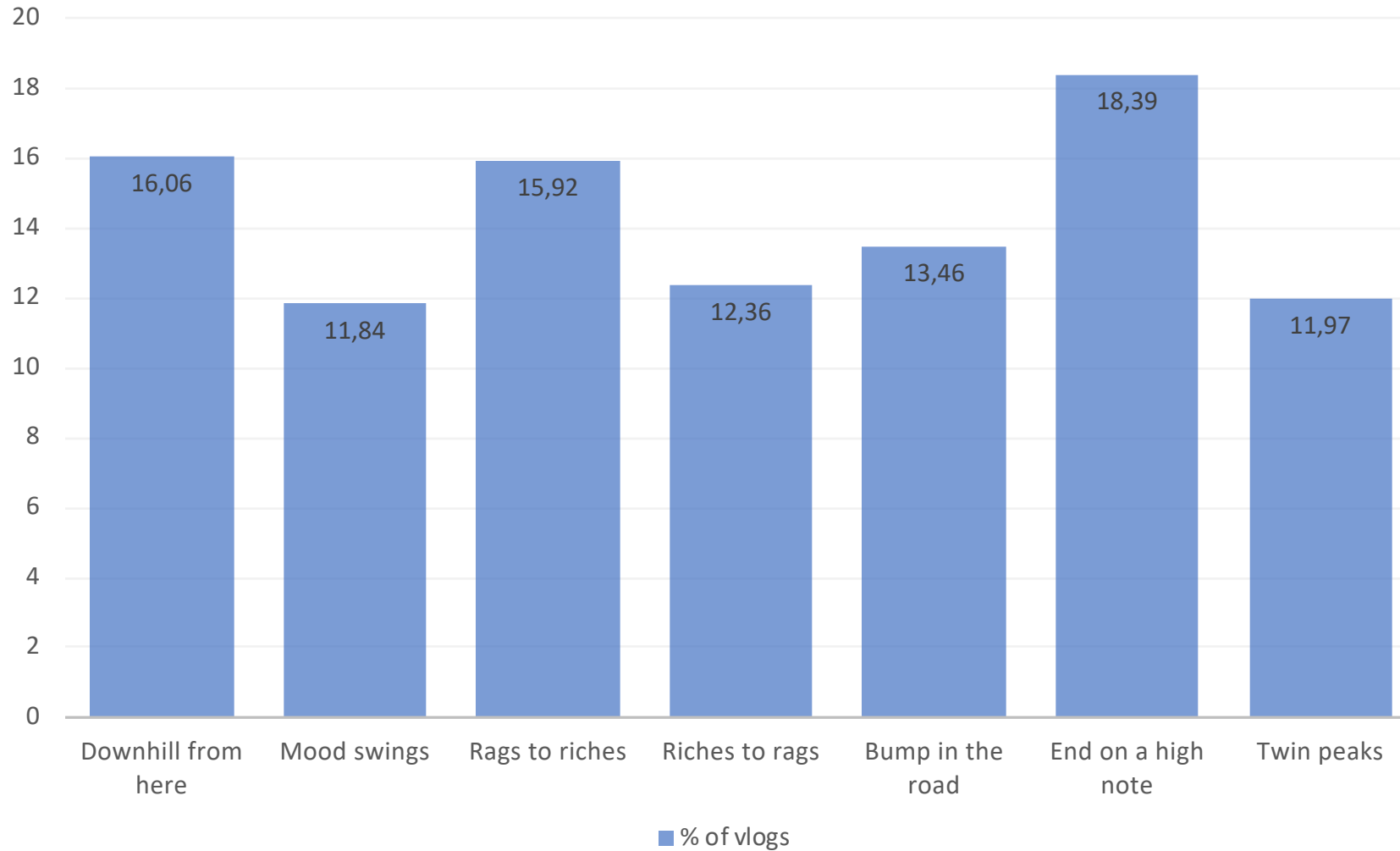
# Results (finding *k)*



Screeplot for k = 1 to k = 30

# Results (finding *k)*
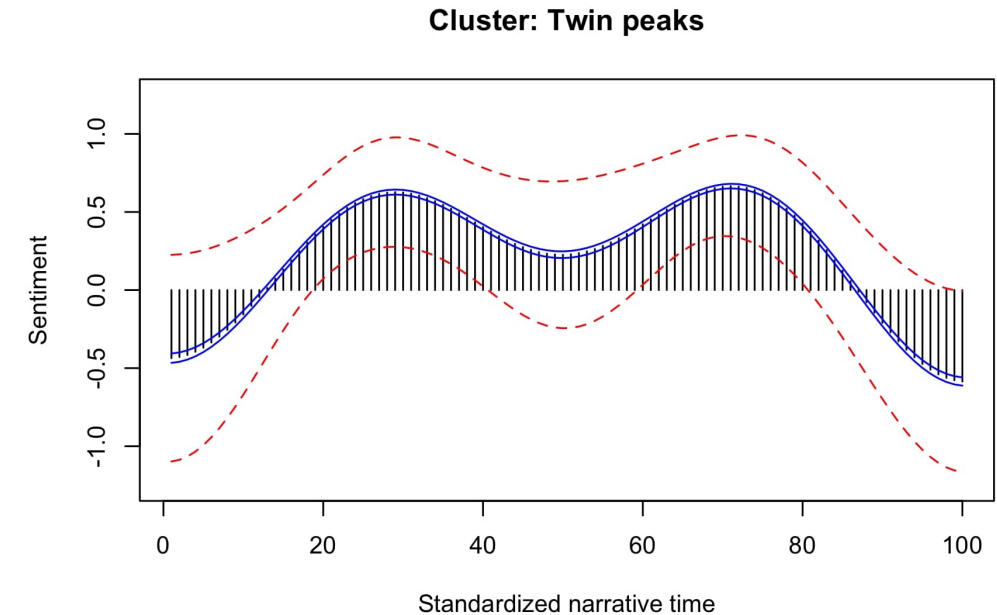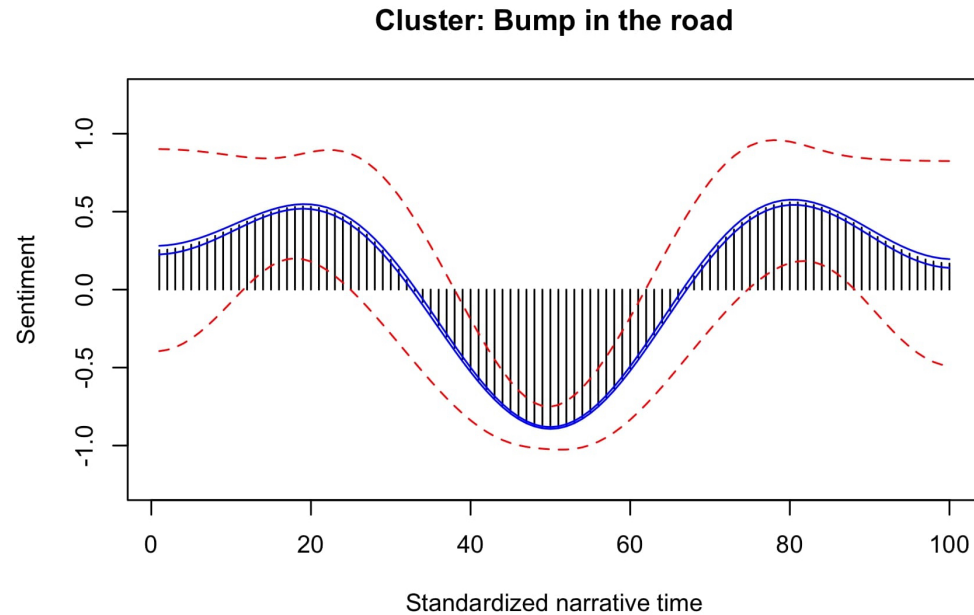


Screeplot for k = 1 to k = 30

Inflection at *k = 7*

# Clusters (*k = 7*)



Video distribution over sentiment clusters

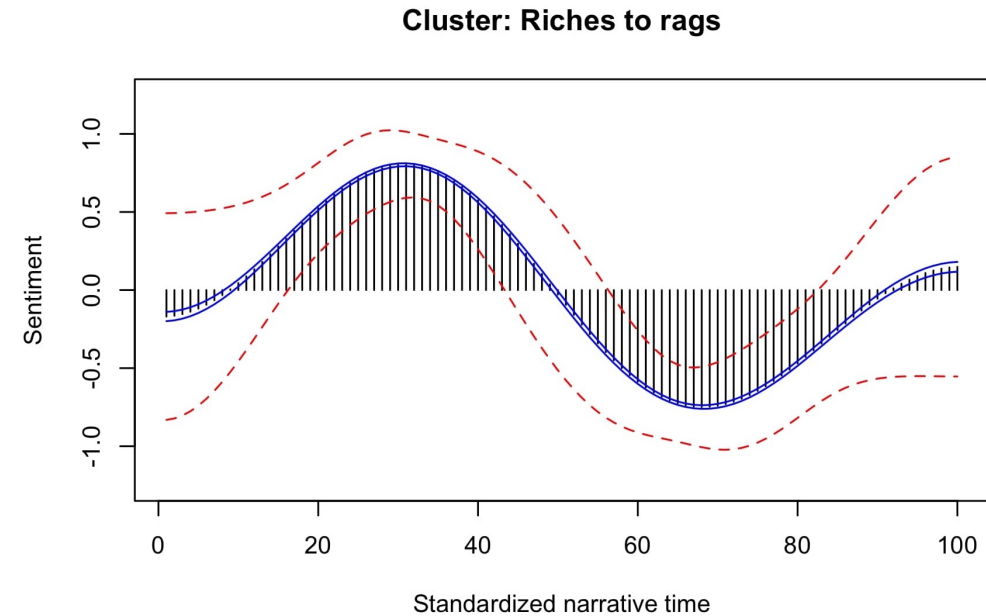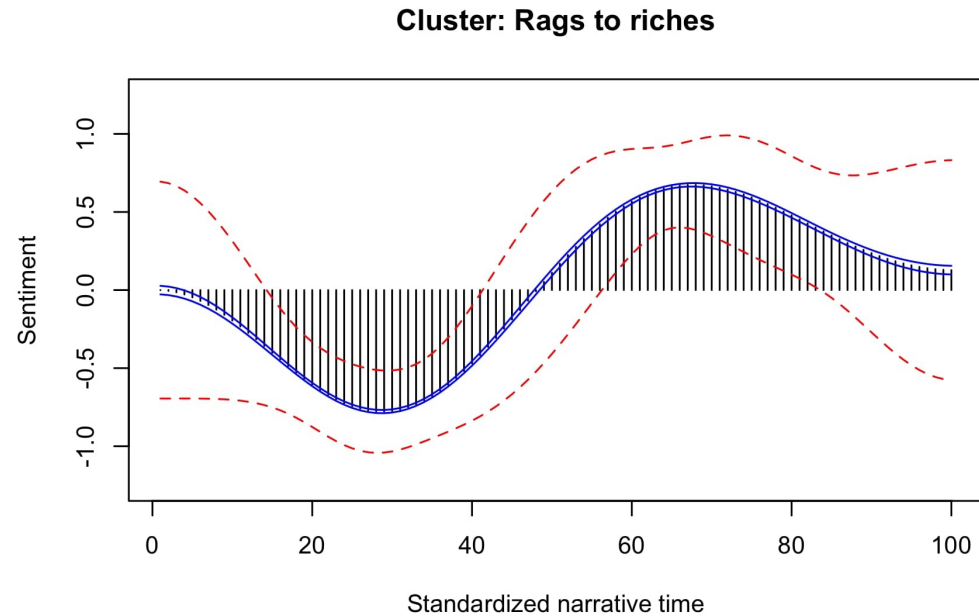# Bump in the road vs. twin peaks

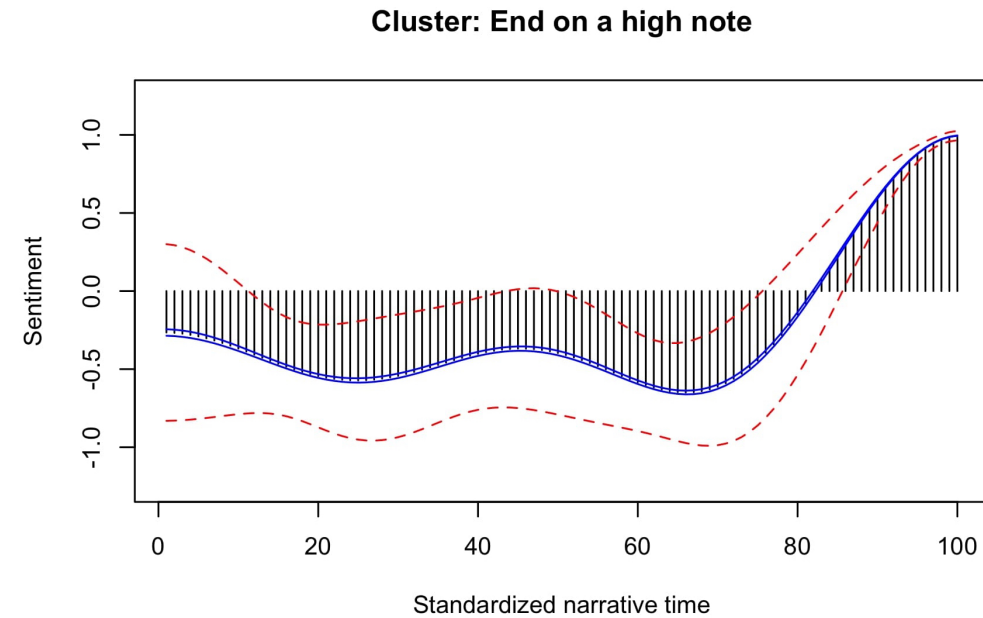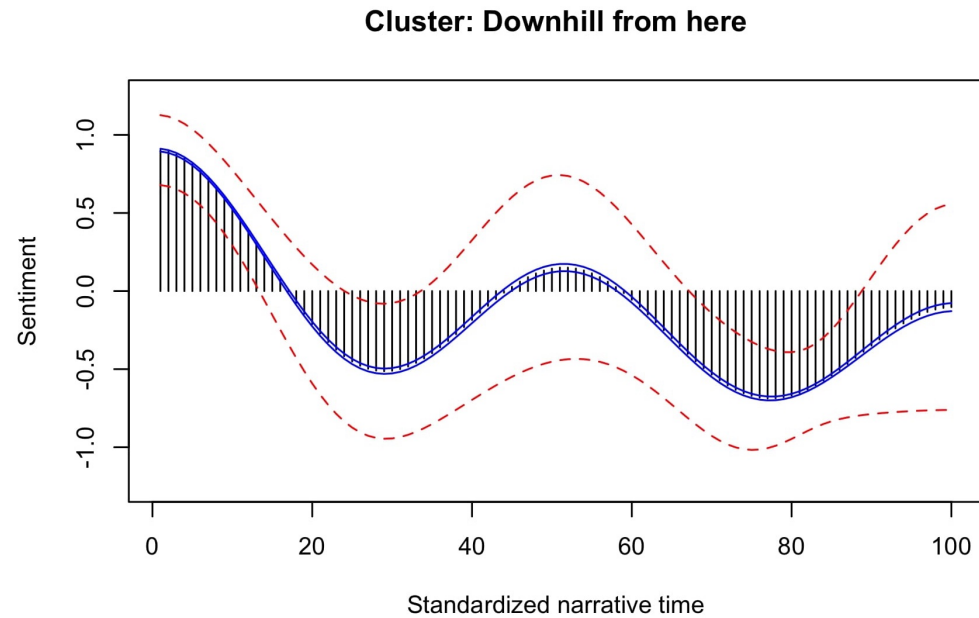**Cluster: Bump in the road**

**Cluster: Twin peaks**



Figures: Average sentiment style shapes. Dotted red lines = +/- 1 SD; blue lines = 99% CI.

# Rags to riches vs. riches to rags



Figures: Average sentiment style shapes. Dotted red lines = +/- 1 SD; blue lines = 99% CI.

# Downhill from here vs. end on a high note

**Cluster: Downhill from here**

**Cluster: End on a high note**

Figures: Average sentiment style shapes. Dotted red lines = +/- 1 SD; blue lines = 99% CI.
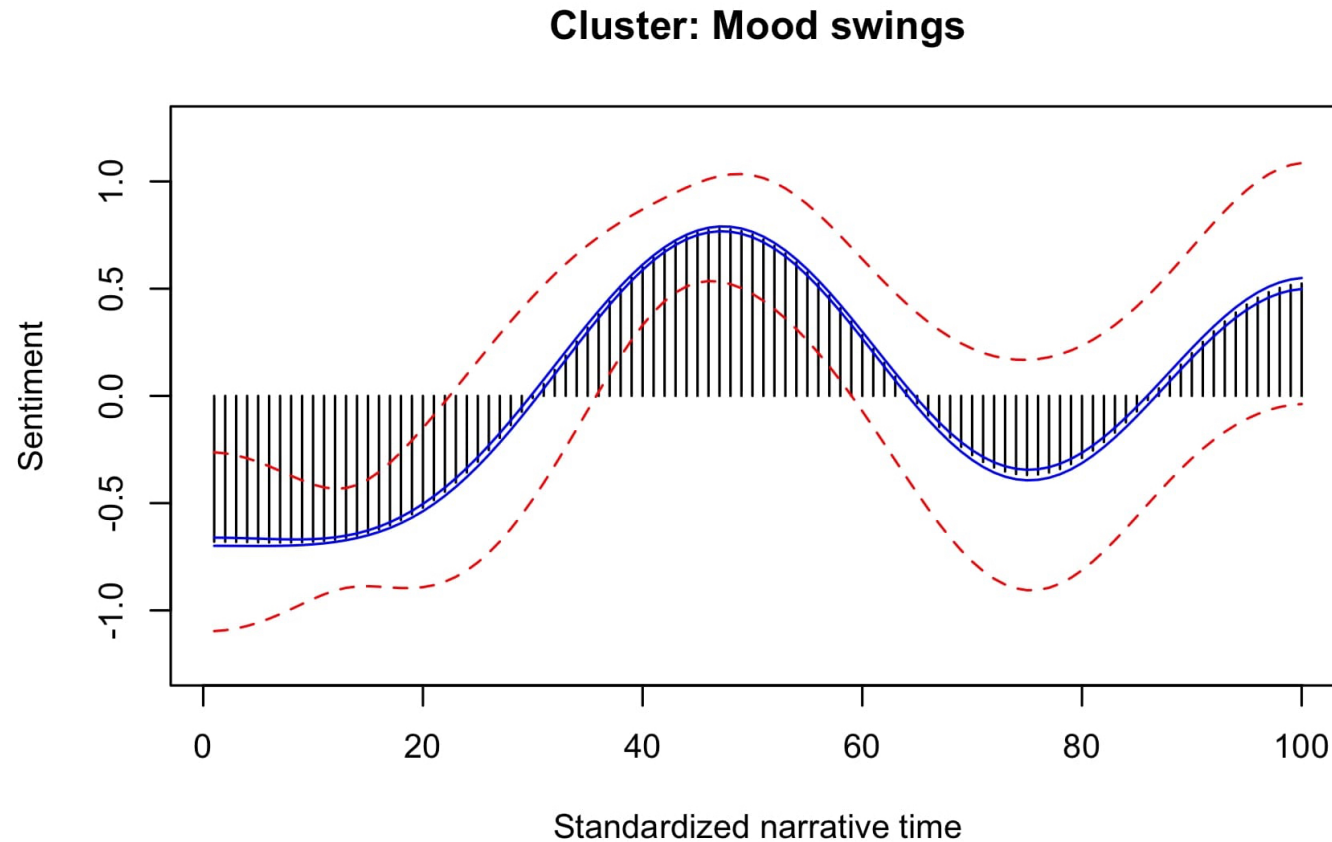
# Mood swings

**Cluster: Mood swings**



Figure: Average sentiment style shapes. Dotted red lines = +/- 1 SD; blue lines = 99% CI.

# Sentiment styles and gender

| | + | - |
|---|---|---|
| Families | twin peaks* | end on a high note* |
| Female | riches to rags* | end on a high note* |
| Male | end-on-a-high-note* | twin peaks*, downhill from here* |

* = significant at $p < .01$

# Limitations

- Automatic transcripts
- Only successful vloggers
- No visual and audio features

# References

- Tyler Rinker. 2018a. lexicon: Lexicon Data. http://github.com/trinker/lexicon

# Resources

- GitHub: https://github.com/ben-aaron188/narrative_structures

- EMNLP publication: https://aclweb.org/anthology/D18-1394