# Sentiment and Econometrics:
## Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

**Samuel Borms**
Université de Neuchâtel
Vrije Universiteit Brussel

August 17, 2020

Public PhD defense

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Motivation

# Motivation

**Sentiment and Econometrics**: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Motivation

# Alternative data

You read the papers. You listen to the radio. You watch TV. Thank you.



Alternative data are "qualitative sentiment data." Information value?!

# Motivation of my work

Focus on **texts**.

**f**(questions, methods, traditional data, textual data) $\rightarrow$ better answers.

But... application-specific textual data transformation is hard.

My thesis attempts to define **f**(). In my version, econometrics meets sentiment meets econometrics.

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Motivation

# User personas – for who is this useful?

The aim is to provide a **gateway** for specialists and non-specialists willing to create and use **textual sentiment** data.
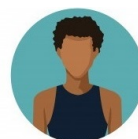
Researcher          Data Scientist          Asset Manager          Macroeconomist

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #1: Formalization of a methodological frame of thought for applied textual (sentiment) analysis

# Contribution #1: Formalization of a methodological frame of thought for applied textual (sentiment) analysis

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #1: Formalization of a methodological frame of thought for applied textual (sentiment) analysis
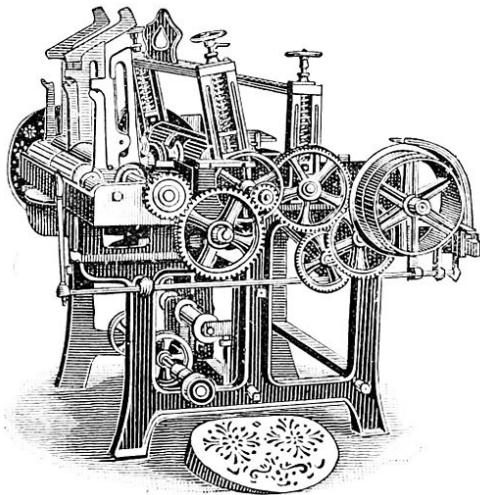
# Generalization of "sentiment"

Step away from the limiting view on sentiment in most literature.

**Definition.** Sentiment is the disposition of an entity toward an entity, expressed via a certain medium.

$\Longrightarrow$ Use the medium... (e.g. press data)
 ... to extract the expressed disposition... (e.g. positivity or bias)
 ... to measure something about one or more entities (e.g. the economy).

Fill in the details along the analysis.

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

└─ Contribution #1: Formalization of a methodological frame of thought for applied textual (sentiment) analysis
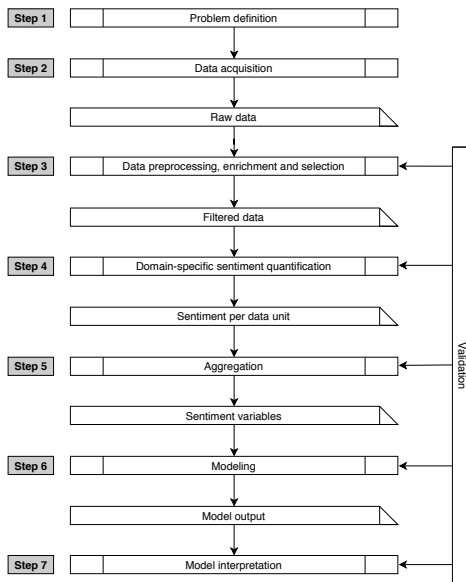
# A framework for solving problems with sentiment data

From problem to data transformation to modeling to concluding and back.

Integrates numerous challenges likely to face.

**sent**iment
+
econ**ometrics**
$\Longrightarrow$
sentometrics

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #1: Formalization of a methodological frame of thought for applied textual (sentiment) analysis

# The sentometrics analysis cycle

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #1: Formalization of a methodological frame of thought for applied textual (sentiment) analysis

# The "joint hypothesis problem" and validation

Main challenge is "joint hypothesis problem"-like. You need to validate both the data transformation and the answer to your research question.

Cyclical validation approach required.

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #2: Formalization of the methodological framework's core into a computational one

**Contribution #2**: Formalization of the methodological framework's core into a computational one

**Sentiment and Econometrics**: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

└─ **Contribution #2**: Formalization of the methodological framework's core into a computational one

# R software package **sentometrics**

Development, release and thorough documentation of open-source R software package **sentometrics**. Free to install and use!



Continuous improvements and additions going forward. See dedicated package website sentometricsresearch.github.io/sentometrics.

**Sentiment and Econometrics**: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

└─ **Contribution #2**: Formalization of the methodological framework's core into a computational one
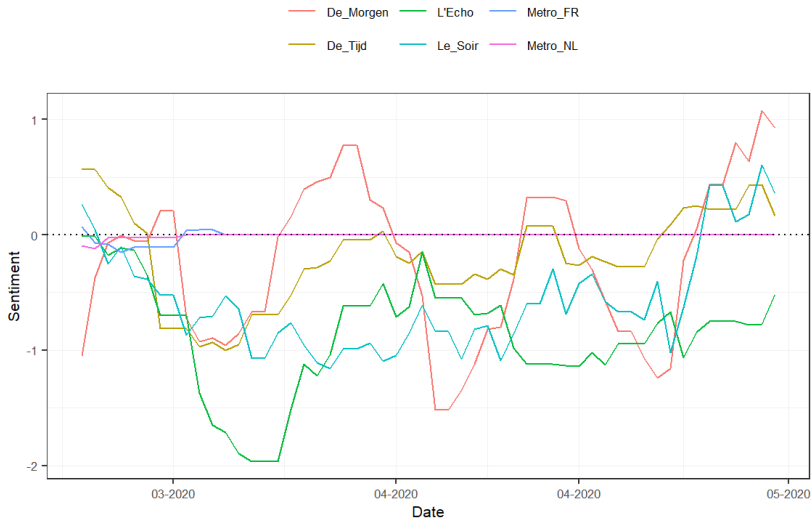
# Functionalities

Maps to methodological framework in that it covers at least one functionality for all steps possible (3–7).

**Unique** when it comes to flexible aggregation into sentiment time series.

| Functionality | Functions | Output |
|---|---|---|
| **1. Corpus management** | | |
| (a) Creation | `sento_corpus()` | *sento_corpus* |
| (b) Manipulation | **quanteda** corpus functions (e.g. `docvars()`, `corpus_sample()`, or `corpus_subset()`), `as.data.frame()`, `as.data.table()`, `as.sento_corpus()` | |
| (c) Features generation | `add_features()` | |
| (d) Summarization | `corpus_summarize()`, `print()` | |
| **2. Sentiment computation** | | |
| (a) Lexicon management | `sento_lexicons()` | *sento_lexicons* |
| (b) Computation | `compute_sentiment()` | *sentiment* |
| (c) Manipulation | `merge()`, `as.sentiment()` | |
| (d) Summarization | `peakdocs()` | |
| **3. Sentiment aggregation** | | |
| (a) Specification | `ctr_agg()` | |
| (b) Aggregation | `sento_measures()`, `aggregate()` | *sento_measures* |
| (c) Manipulation | `subset()`, `merge()`, `diff()`, `scale()`, `as.data.frame()`, `as.data.table()`, `measures_fill()`, `measures_update()` | |
| (d) Visualization | `plot()` | |
| (e) Summarization | `summary()`, `peakdates()`, `print()`, `nobs()`, `nmeasures()`, `get_dimensions()`, `get_dates()` | |
| **4. Modeling** | | |
| (a) Specification | `ctr_model()` | |
| (b) Estimation | `sento_model()` | *sento_model*, *sento_modelIter* |
| (c) Prediction | `predict()` | |
| (d) Diagnostics | `summary()`, `print()`, `get_loss_data()`, `attributions()` | *attributions* |
| (e) Visualization | `plot()` | |

**Sentiment and Econometrics**: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

└─**Contribution #2**: Formalization of the methodological framework's core into a computational one

# Litmus test – example sentiment time series variables

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #3: Structured news-based measurement of firm-level sustainability and of economic uncertainty

# Contribution #3: Structured news-based measurement of firm-level sustainability and of economic uncertainty

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #3: Structured news-based measurement of firm-level sustainability and of economic uncertainty

# Application #1: News-based firm-level sustainability

Increasing interest to invest in companies who do well environmentally, socially and governance wise (ESG).

**Sustainable asset managers** use external ESG ratings and in-house research to screen the investment universe.

I add news to the mix through creation of daily firm-level indices capturing frequency and sentiment of news reporting about ESG issues.

# Empirical analysis

291 European stocks, Dutch news from Belga, 1999–2018, Sustainalytics.

Use of specific keywords to track the news relevant to ESG **coverage**.
More negative **sentiment** word list ("controversies").

Querying, cleaning, selection, aggregation, validation.

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #3: Structured news-based measurement of firm-level sustainability and of economic uncertainty

# Keywords

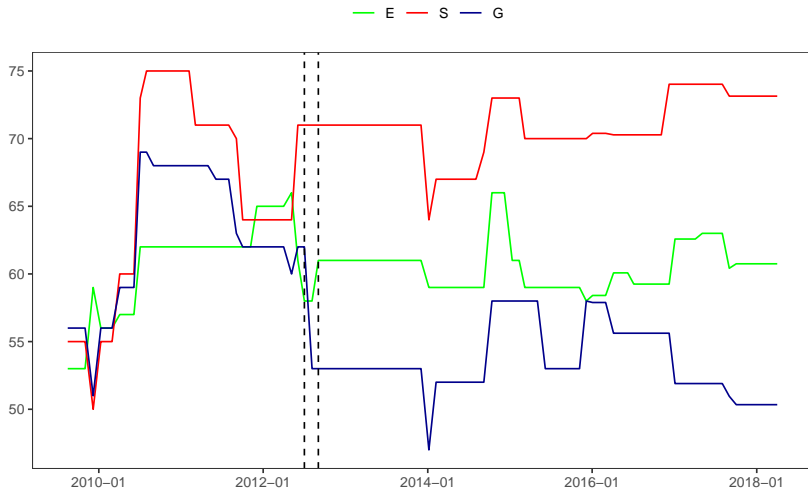A **semi-supervised** approach: expertise + models. Man & machine.

Expert: pick some important words (klimaat, mobiliteit, ecologie, etc.).

Model: from >100000 words, tell me which words are semantically related based on estimated **word embeddings**.
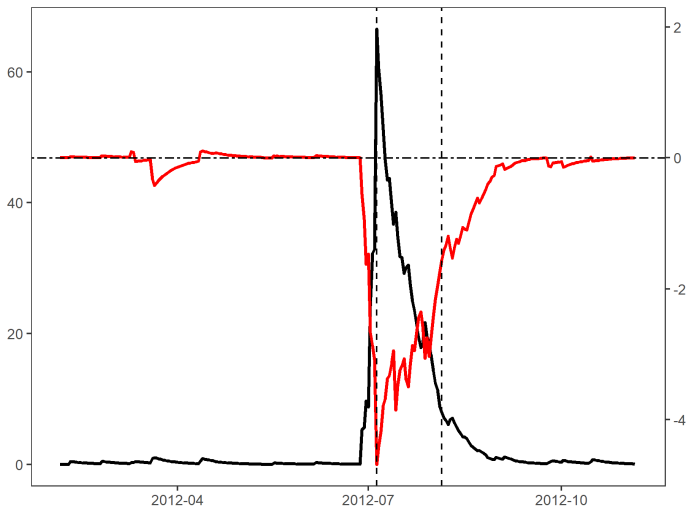
Expert: check if the most related words are useful.

Algorithm: follow up the news that use these keywords.

# Barclays Sustainalytics ratings (monthly)

# Barclays news coverage & sentiment during LIBOR scandal



Short-term news **signals**. View of the risk-averse doctor analyzing patients.

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #3: Structured news-based measurement of firm-level sustainability and of economic uncertainty

# Stock and sector screening

Monthly rebalanced portfolios based on the news-based indicators perform at least as well as portfolios based on external ESG ratings.

News coverage indicators more informative than pure sentiment ones.
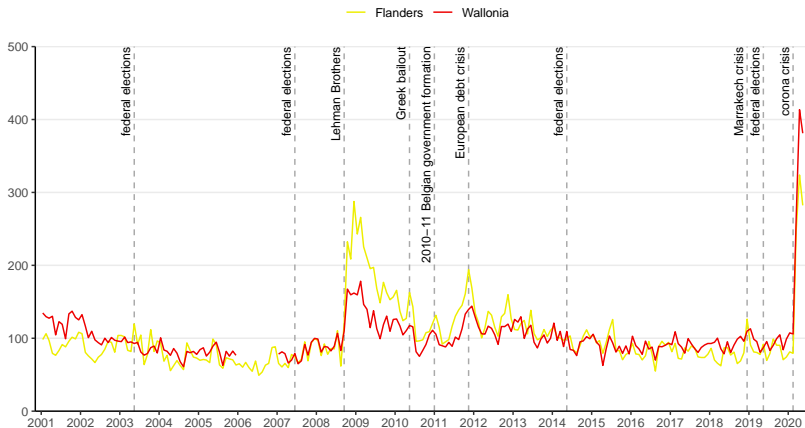
Sector rotation most promising.

*Possible extensions*: international news, event study, factor portfolios.

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #3: Structured news-based measurement of firm-level sustainability and of economic uncertainty
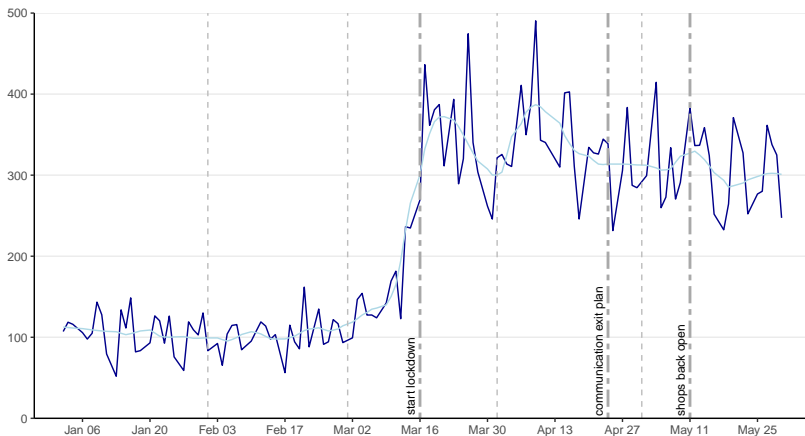
# Application #2: An EPU index for Belgium

Application to the case of Belgium of common methodology to measure economic policy uncertainty (EPU) from news articles.

Same word embedding approach to keywords definition.

**Sentiment and Econometrics**: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

└─ **Contribution #3**: Structured news-based measurement of firm-level sustainability and of economic uncertainty

# Monthly evolution of news-based EPU in Belgium

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #3: Structured news-based measurement of firm-level sustainability and of economic uncertainty

# Daily zoom-in during 2020

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #3: Structured news-based measurement of firm-level sustainability and of economic uncertainty

# Explaining peaks

Additional validation in the form of automated qualitative "news reader."

1    bedrijven, miljoen, coronacrisis, bedrijf, miljard, maand, banken, werknemers, week, België

2    week, N-VA, land, mei, weken, tijd, Veiligheidsraad, coronacrisis, CD&V, leven

3    virus, land, landen, wereld, aantal, lockdown, China, Trump, leven, coronavirus

4    landen, miljard, Italië, geld, Nederland, Europa, bedrijven, land, coronacrisis, EU

5    miljoen, Brussels Airlines, coronacrisis, bedrijven, stad, vraag, Lufthansa, weken, geld, mei

5 clusters of news in April 2020, all related to COVID-19.

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Contribution #3: Structured news-based measurement of firm-level sustainability and of economic uncertainty

# www.policyuncertainty.com



But what in August 2020? And so on...

**Sentiment and Econometrics**: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

└─ Conclusion

# Conclusion

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Conclusion

# Contributions wide and large

A systematic approach to use textual data in applied research.

A computational toolbox allowing to do so quickly and efficiently.

A catalyst to more effectively use textual data.

An inspiration to formalize similar frameworks for audio and video data.

A structured application to following ESG-related news.

A monthly index to track EPU in Belgium.

**Sentiment and Econometrics**: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Conclusion

# Exciting future research

Better **validation** tools.

More **applications** based on the framework.

Integrated **theoretical developments** jointly accounting for various steps.

**Intraday** textual sentiment analysis.

**Multimodal** sentiment analysis.

Wouldn't it be nice if we could collaborate cross-disciplinary?

Sentiment and Econometrics: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Conclusion

# github.com/SentometricsResearch (in progress)

# PhD papers

▶ Algaba, Ardia, Bluteau Borms & Boudt (2020). "**Econometrics meets sentiment: An overview of methodology and applications**". The Journal of Economic Surveys 34 (3), 512-547.

▶ Ardia, Bluteau, Borms & Boudt (2020). "**The R package *sentometrics* to compute, aggregate and predict with textual sentiment**". The Journal of Statistical Software, forthcoming.

▶ Borms, Boudt, Van Holle & Willems (2020). "**Semi-supervised text mining for monitoring the news about the ESG performance of companies**". Data Science for Economics and Finance: Methodologies and Applications (Springer), forthcoming.

▶ Algaba, Borms, Boudt & Van Pelt (2020). "**The Economic Policy Uncertainty index for Flanders, Wallonia and Belgium**". Bank- en Financiewezen digitaal 2020/6.

**Sentiment and Econometrics**: Toward A Unified Framework of Textual Sentiment Analysis for Economic and Financial Applications

Conclusion

# Thanks

Many thanks for your attention!

Looking forward to taking questions.