# Extracting Socioeconomic Patterns from the News: Modelling Text and Outlet Importance Jointly

**Vasileios Lampos**[β]
v.lampos@ucl.ac.uk

**Daniel Preoţiuc-Pietro**[δ]
danielpr@gmail.com

**Sina Samangooei**[σ]
sinjax@gmail.com

**Douwe Gelling**[γ]
d.gelling@sheffield.ac.uk

**Trevor Cohn**[τ]
t.cohn@unimelb.edu.au

UCL[β]   The University Of Sheffield.[δ,γ]   University of Southampton[σ]   Trend Miner   THE UNIVERSITY OF MELBOURNE[τ]

$$f : X \rightarrow Y$$

$$\{(\text{📄}, \text{OUTLET})\} \quad ESI^*$$

2006 -> 2013

prediction — outlet x word daily frequencies — bias term

$$y = oXw + \beta$$

outlet weights (o)



n-gram weights (w)



## Forecast



Linear regression: 9.253 (9.89%)
Bilinear reression: 8.209 (8.77%)

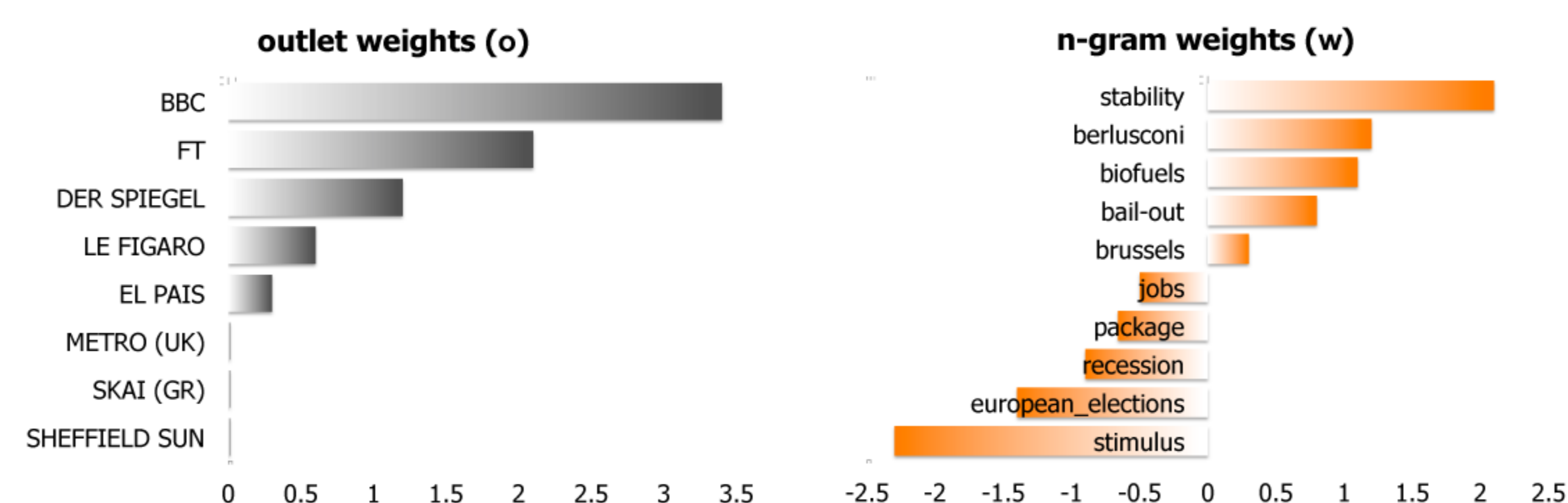RMSE (error rate %) in a forecasting setup on 30 data points

## Data

- news summaries written in English
- extracted from the **Open Europe Think Tank**
- daily aggregat1ion of news about the EU
- focused on current affairs, politics, economy
- each summary is a few paragraphs long
- each summary is attached to ≥1 news outlets
- **1913 days**, 94 months
- avg. 14 news/summary
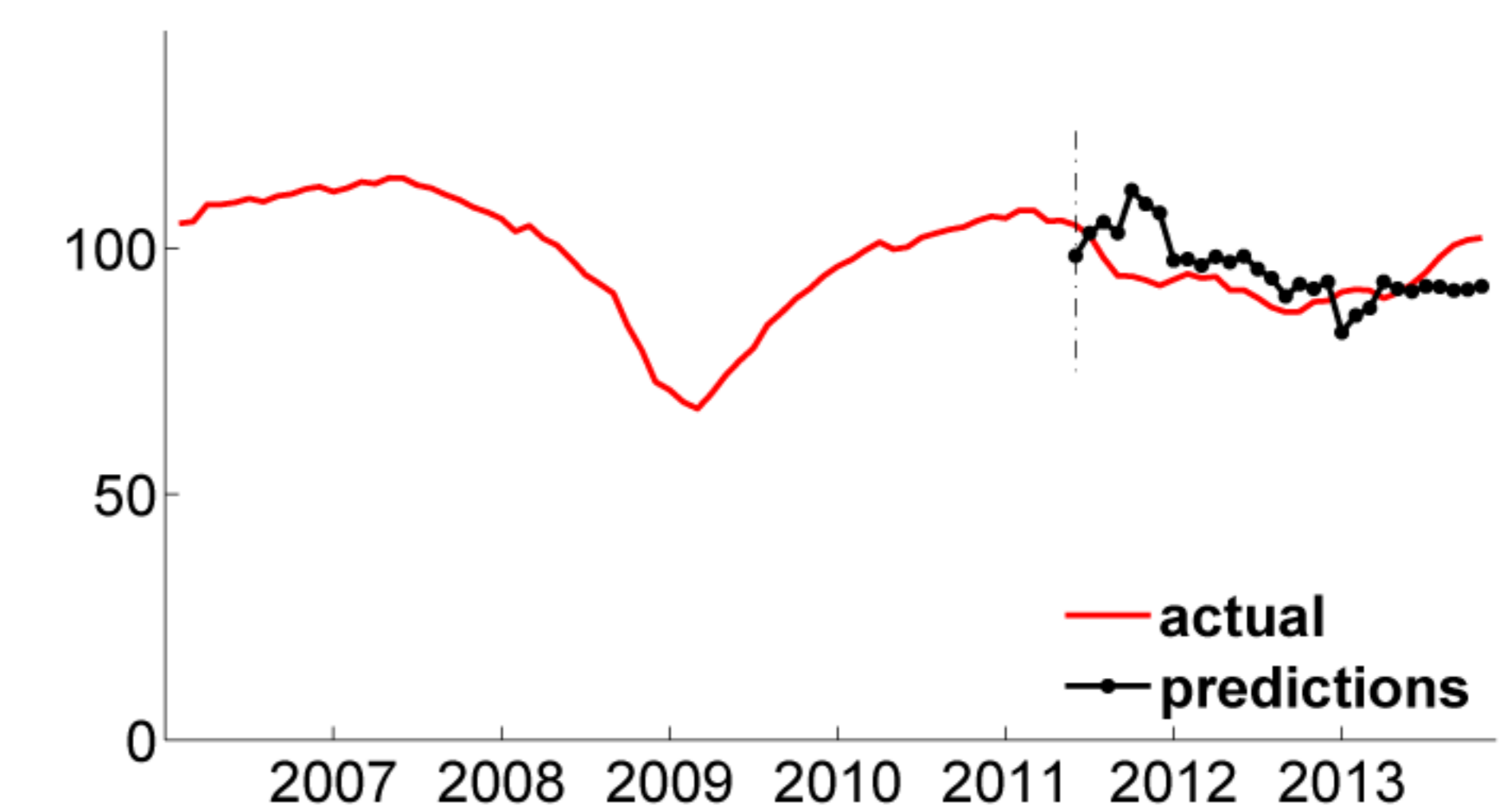- features: **8413 unigrams + 19045 bigrams**
- **435 news outlets**

**\*ESI = Economic sentiment indicator**
- composite indicator often seen as an early predictor for future economic developments
- consists of five confidence indicators with different weights: industrial (40%), services (30%), consumer (20%), construction (5%) and retail trade (5%)
- any socioeconomic indicator can be used in this framework



Frequency    Weight    Polarity
Word         a → **a**   +      -
Outlet                    Yes    Yes