# Gavagai

# Text Analysis for Big Data
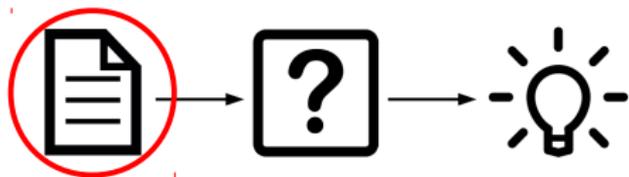
Magnus Sahlgren

# Text analysis

# Text analysis

# Text analysis

Size

Style
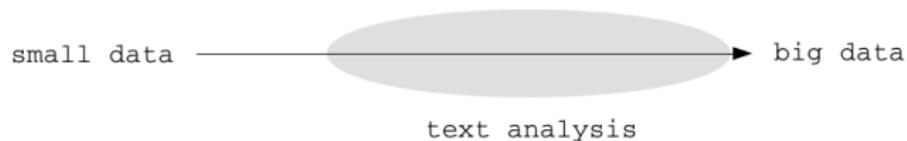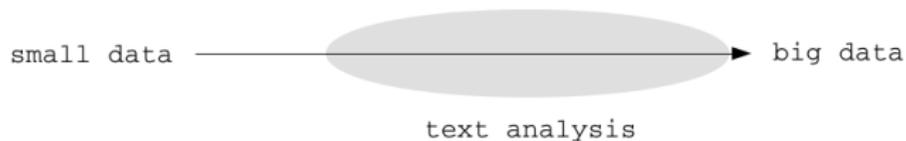(editorial vs social)

Language
(there are other languages than English out there)

Gavagai

# Text analysis

Size

Style
(editorial vs social)

Language
(there are other languages than English out there)

Gavagai

# Text analysis

## Size

small data ——————⟶ big data

text analysis

## Style
(editorial vs social)

## Language
(there are other languages than English out there)

# Text analysis

## Size

small data ────────⟶ big data

text analysis

## Style

(editorial vs social)
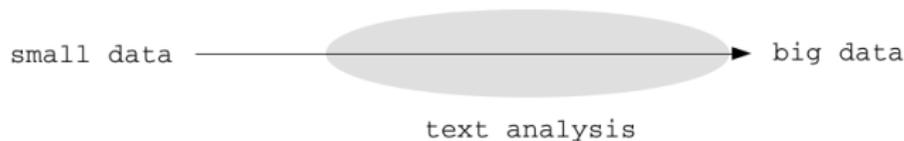
Language

(there are other languages than English out there)

Gavagai

# Text analysis

## Size

small data ———————➤ big data

text analysis

## Style
(editorial vs social)

Language
(there are other languages than English out there)

Gavagai

# Text analysis

### Size

small data ——————→ big data

text analysis

### Style
(editorial vs social)

### Language
(there are other languages than English out there!)

Gavagai

# Text analysis

### Size

small data ——————→ big data

text analysis

### Style
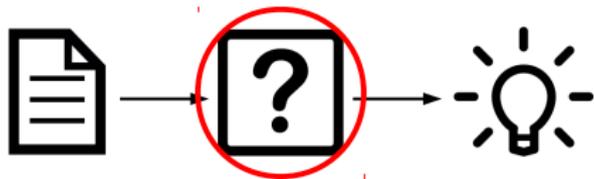(editorial vs social)

### Language
(there are other languages than English out there!)

Gavagai

# Text analysis

# Text analysis

Knowledge-based
(use resources like Wikipedia)

Supervised machine learning
(use annotated data)

Unsupervised machine learning
(use unstructured data)

Gavagai

# Text analysis

### Knowledge-based
(use resources like Wikipedia)

Supervised machine learning
(use annotated data)

Unsupervised machine learning
(use unstructured data)

Gavagai

# Text analysis

### Knowledge-based
(use resources like Wikipedia)

### Supervised machine learning
(use annotated data)

### Unsupervised machine learning
(use unstructured data)

Gavagai

# Text analysis

Knowledge-based
(use resources like Wikipedia)

Supervised machine learning
(use annotated data)

Unsupervised machine learning
(use unstructured data)

Gavagai

# Text analysis

**Knowledge-based**
(use resources like Wikipedia)

**Supervised machine learning**
(use annotated data)

Unsupervised machine learning
(use unstructured data)

Gavagai

# Text analysis

Knowledge-based
(use resources like Wikipedia)

Supervised machine learning
(use annotated data)

Unsupervised machine learning
(use unstructured data)

Gavagai

# Text analysis

Knowledge-based
(use resources like Wikipedia)

Supervised machine learning
(use annotated data)

Unsupervised machine learning
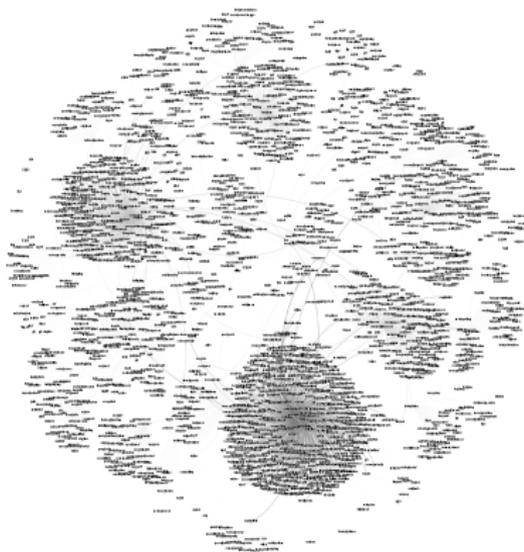(use unstructured data)

Gavagai

# Semantic memories

# Semantic memories

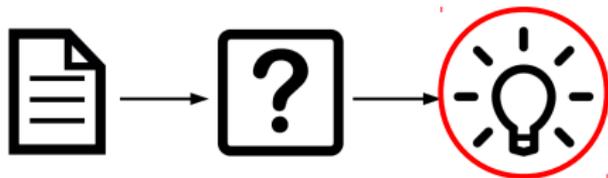(systems that learn language by reading large amounts of text)

Gavagai

# Semantic memories

(systems that learn language by reading large amounts of text)



Gavagai

# Text analysis

# Text analysis

Identify and extract items
(e.g. entities and events)

Find relations
(e.g. synonyms and associations)

Gavagai

Identify and extract items
(e.g. entities and events)

Find relations
(e.g. synonyms and associations)

Gavagai

# Text analysis

Identify and extract items
(e.g. entities and events)
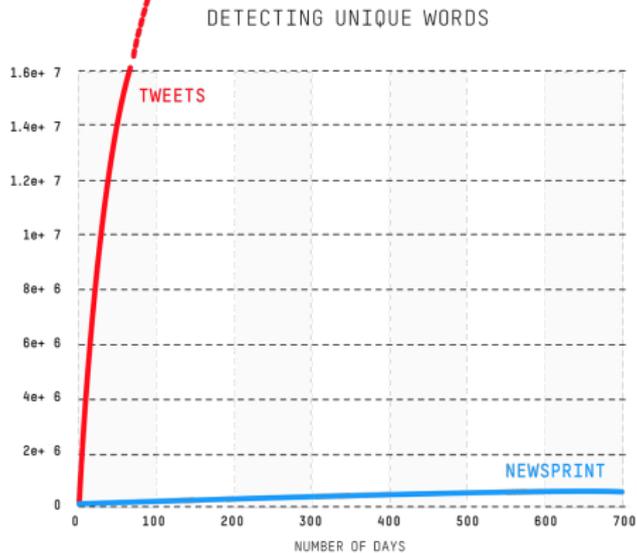
Find relations
(e.g. synonyms and associations)

Gavagai

# Text analysis

### Identify and extract items
(e.g. entities and events)

### Find relations
(e.g. synonyms and associations)

Gavagai

# Text analysis

### Identify and extract items
(e.g. entities and events)

### Find relations
(e.g. synonyms and associations)

Gavagai

# Find relations

# Find relations



DETECTING UNIQUE WORDS

TWEETS

NEWSPRINT

Gavagai

# Find relations (`lexicon.gavagai.se`)

GAVAGAI LIVING LEXICON

Type a word, choose language and hit Enter

economy | English | Enter

economy is the **806** most frequent term in our **English** lexicon, which puts it at the top **0.039%** of the vocabulary that contains some **2,399,006** different terms.

Below are lists of terms that our semantic memories find related to the term you entered. Click on a term to see examples of how they are used.

LIKE WHAT YOU SEE?
Try our **Chrome Extension** to be able to look up words that you find when you are browsing the web.

Want to implement our word knowledge in your own applications? **Sign up** for a free trial of our API.

**SIMILARLY SPELLED**
These words look similar to economy

econom
economy's
econo
economia
economie
economic
econam
ectomy
ecology
ekonomi

**LEFT SIDE NEIGHBOURS**
These terms are frequently used before economy

eurozone
slowing
bitcoins
the greek
russia's
the chinese
greece's

**SEMANTICALLY SIMILAR**
These terms are used with the same other terms as economy

(india's)
manufacturing sector              (contracted)
gross domestic product
defence budget
japan's economy
foreign exchange reserves
industrial sector

(the greek, greece's)
banking sector
financial system
banking system
economic output
financial sector
economic woes
gold reserves
economic crisis
banking industry
ruling elite
debt crisis
stock markets
real economy
pension system
finance ministry

(tunisia's)                        (government-held)
agricultural sector
tourism sector
agriculture sector
energy mix
tourism industry
second-largest
economic base
nuclear industry

(slowing, sluggish)
domestic economy
local economy
global economy
economic growth

property market                    (personal finance)
housing market

**RIGHT SIDE NEIGHBOURS**
These terms are frequently used after economy

contracted
is it you

**N-GRAMS**
economy is used as part of these terms

china's economy
the global economy
fuel economy
world's second-largest economy
largest economy
global economy
the world economy
chinese economy
sharing economy
indian economy
digital economy
world's second-biggest economy
domestic economy
political economy
canada's economy
slowing chinese economy
economy class
japan's economy
local economy
india's economy
world economy
premium economy
circular economy
market economy
rural economy
national economy
second-largest economy
british economy
asia's third-largest economy
new economy

**ASSOCIATIONS**
These terms are currently used in the same texts as economy

moneyweek
headlines yahoo
the global economy
schwedule
economies
euro zone
ftse 100

recession
canada's economy
consecutive quarters
statistics canada
canada's

world's second-largest economy

53.4
40.7
fastest pace

stock indexes

china's economy

stephen harper
canadians

world's most populous
spurred

shrank

dethroned

growth model

hard landing

market forces

lurching

crunch

Gavagai

# Text analysis

Identify and extract items
(e.g. entities and events)

Find relations
(e.g. synonyms and associations)

Compress and refine the information
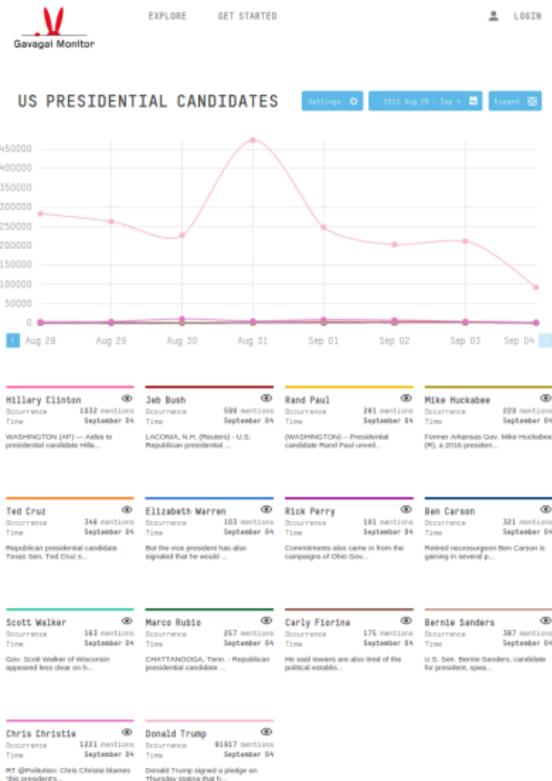(e.g. summarization and topic detection)

Gavagai

# Text analysis

Identify and extract items
(e.g. entities and events)

Find relations
(e.g. synonyms and associations)

Compress and refine the information
(e.g. summarization and topic detection)

Gavagai

# Text analysis

Identify and extract items
(e.g. entities and events)

Find relations
(e.g. synonyms and associations)

Compress and refine the information
(e.g. summarization and topic detection)

Gavagai

# Compress and refine

Summarization and topic detection

# Compress and refine (`monitor.gavagai.se`)

Summarization and topic detection



Gavagai

# Compress and refine (`monitor.gavagai.se`)

Summarization and topic detection



Gavagai

# Compress and refine (`monitor.gavagai.se`)

## Summarization and topic detection



Gavagai

# Text analysis

Insights

Identify and extract items
(e.g. entities and events)

Find relations
(e.g. synonyms and associations)

Compress and refine the information
(e.g. summarization and topic detection)

Measure things
(e.g. attitudes and opinions)

Gavagai

# Text analysis

Identify and extract items
(e.g. entities and events)

Find relations
(e.g. synonyms and associations)

Compress and refine the information
(e.g. summarization and topic detection)

Measure things
(e.g. attitudes and opinions)

Gavagai

# Text analysis

Identify and extract items
(e.g. entities and events)

Find relations
(e.g. synonyms and associations)

Compress and refine the information
(e.g. summarization and topic detection)

Measure things
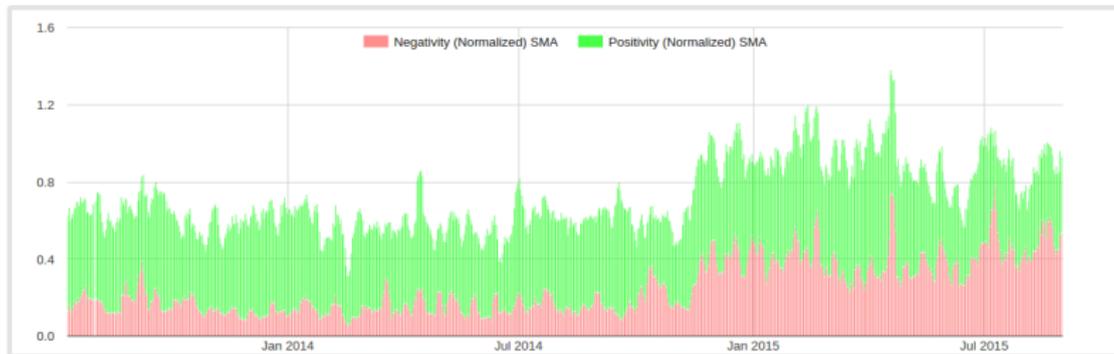(e.g. attitudes and opinions)

Gavagai

# Measure
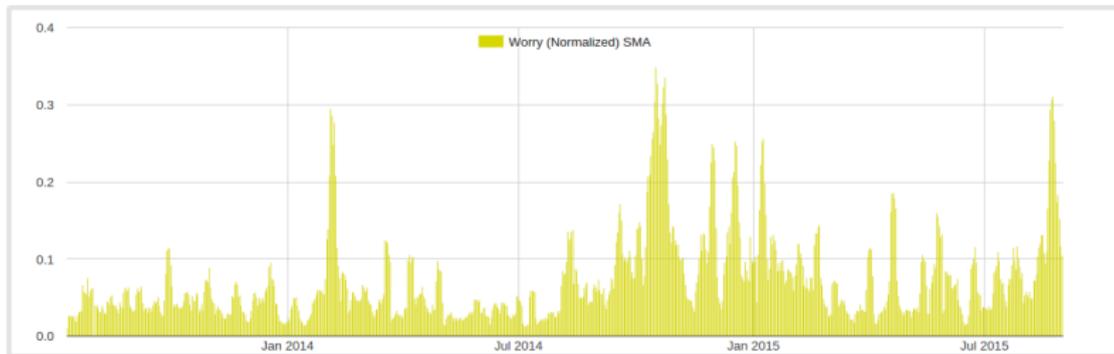
Sentiment analysis

Gavagai

# Measure
## Sentiment analysis



Positivity vs negativity wrt the global economy in English online media

Gavagai
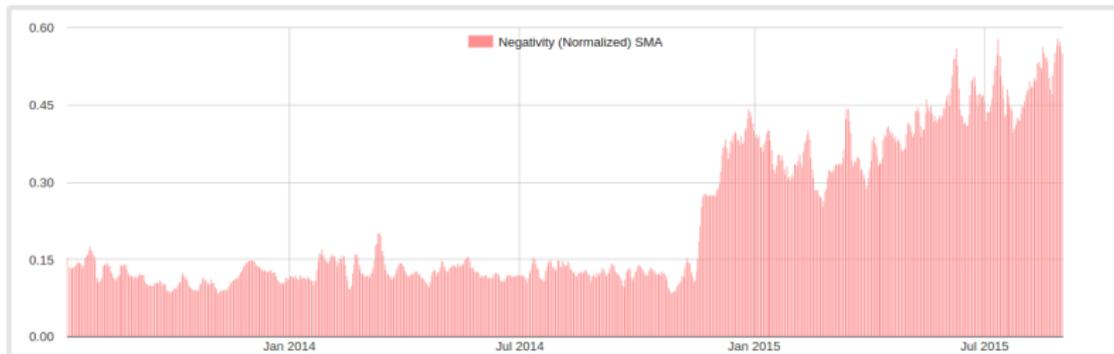
# Measure

Worry wrt the global economy in English online media

Gavagai

# Measure
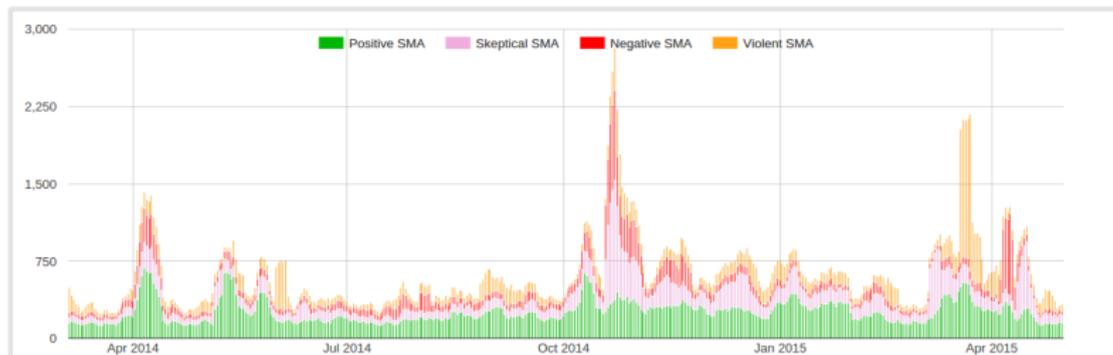
## Sentiment analysis



Negativity towards China in English online media

Gavagai

# Measure

Attitude towards Sweden in Russian online media

Gavagai

# Measure

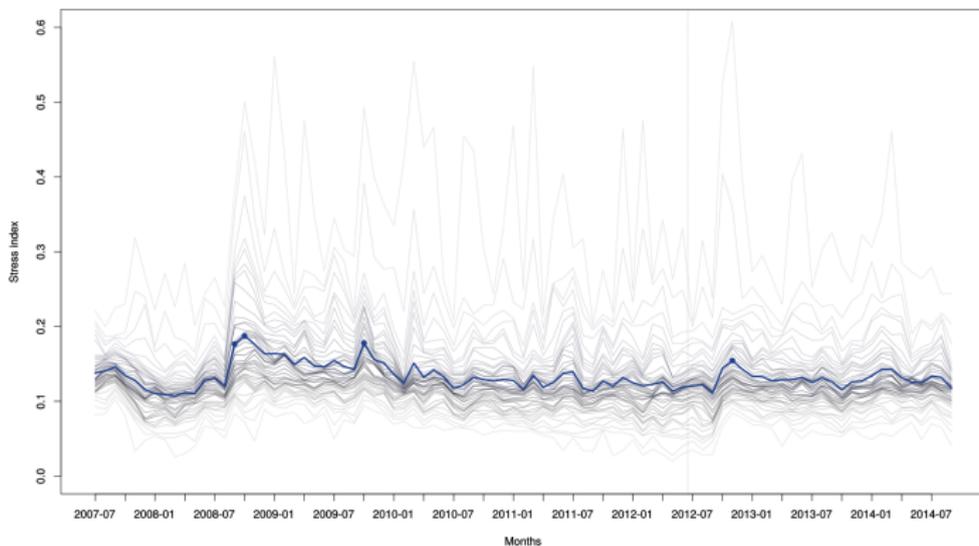Predict

Gavagai

# Measure

Rönnqvist & Sarlin (2015): Detect & Describe: deep learning of bank stress in the news

Gavagai