# TOWARDS SOCIAL MEDIA MINING: TWITTEROBSERVATORY

*Inna Novalija, Miha Papler, Dunja Mladenić*

Artificial Intelligence Laboratory
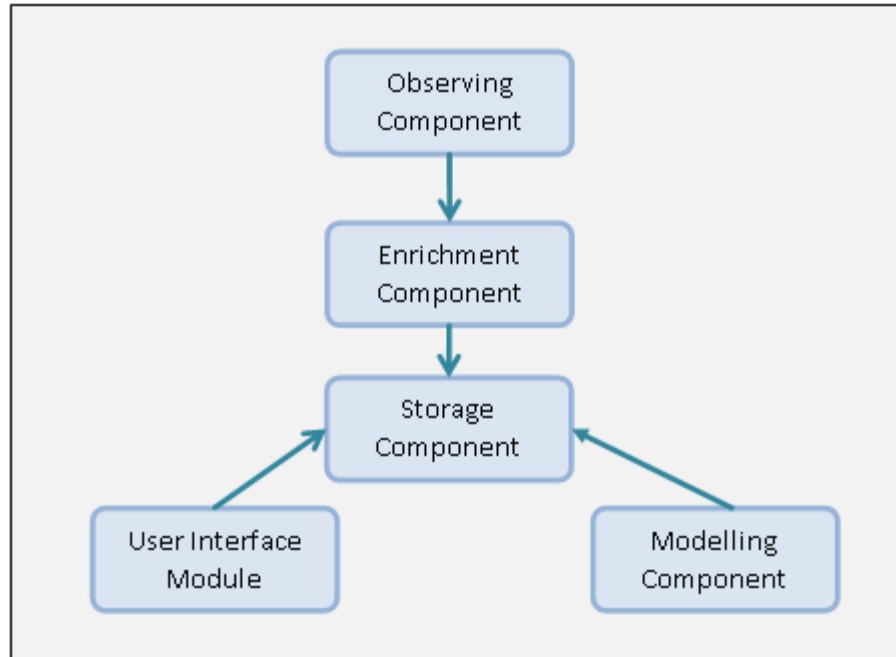
Jožef Stefan Institute

ailab.ijs.si

# Introduction

- **Goal** of the **Social Media Mining** is data mining of **content** streams **produced by people** through **interaction** via Internet based applications.

- **Social media mining** is usually associated with **noisy, distributed, unstructured and dynamic data**, as well as with informal text processing.

- In this research we introduce a novel **Social Media Mining Pipeline** and **TwitterObservatory** tool for:
    - **observing**,
    - **enriching**,
    - **storing**,
    - **analyzing** and
    - **presenting** information

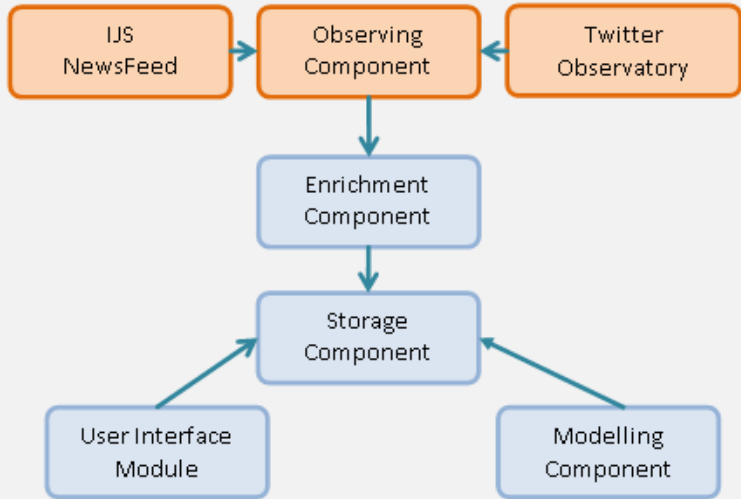obtained from social media and in particular, from Twitter.
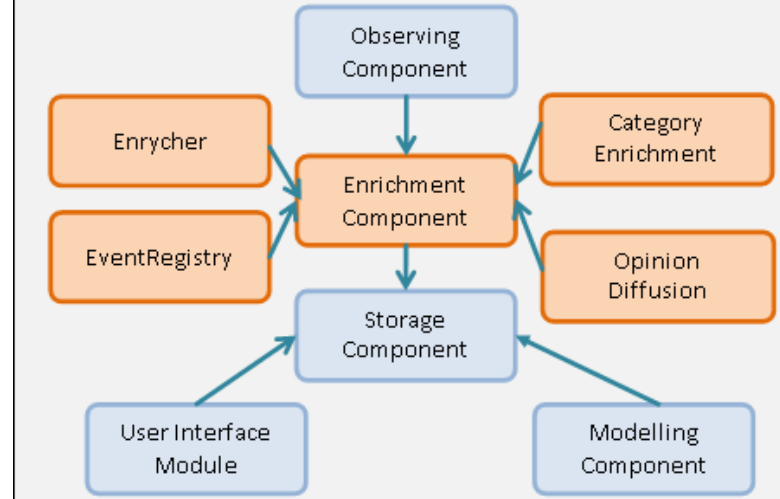
# Social Media Mining Pipeline
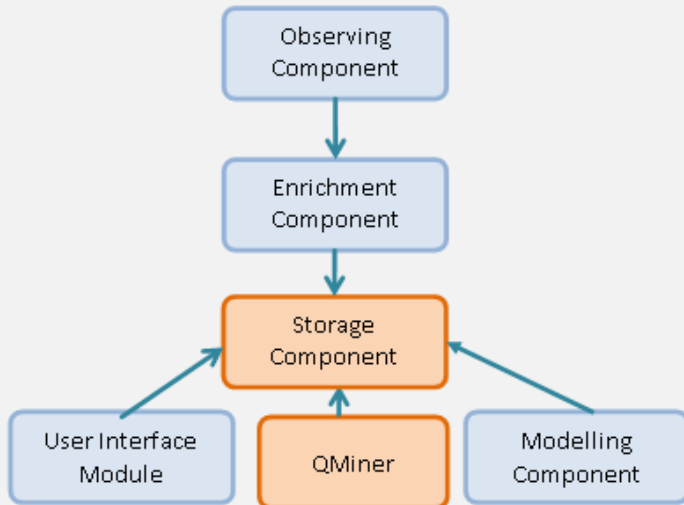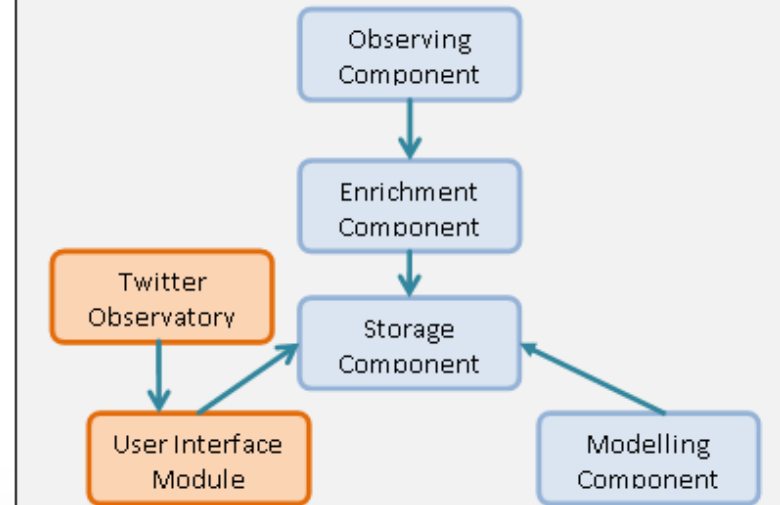
# Social Media Mining Pipeline

# Observing

o **OBSERVING SOCIAL MEDIA by LOCATION:**
- Geo coordinates from United Kingdom
- 10 largest cities (by population)

| Location | Number of Tweets (archived) |
|---|---|
| United Kingdom | 31 GB |

o **OBSERVING SOCIAL MEDIA by KEYWORDS:**
- 400 most common words from Wikipedia,

| Keywords | Number of Tweets (archived) |
|---|---|
| Common words | 500 GB |

# Enriching

o **DATA ENRICHMENT with ENRYCHER for TwitterObservatory:**

- **Enrycher** is a service-oriented system that aims to **shallow and deep text processing** functionalities.
- **Enrycher** processing **functionalities** include:
  - o **Topic and keyword detection**
  - o **named entity extraction**: names of people, locations and organizations, dates, percentages and money amounts
  - o **sentiment** enrichment (English language)
  - o Etc.

o **CATEGORY ENRICHMENT with XLING for TwitterObservatory:**

  - o **Cross-lingual** DMOZ categorization

# Storing

o **QMINER** as storing and analytics platform **for TwitterObservatory**.

o **QMiner** is a **data analytics platform** for **streams of structured and unstructured** data that at the same time contains a number of techniques for **supervised, unsupervised and active learning** on streams of data.

```
[
{
        "name" : "SocialMediaInput",
        "fields" : [
                { "name" : "URI", "type" : "string", "primary" : true },

                { "name" : "Language", "type" : "string", "codebook" : true },
                { "name" : "DateTime", "type" : "datetime" },
                { "name" : "Title", "type" : "string", "store" : "cache" },
                { "name" : "Body", "type" : "string", "store" : "cache" },

                { "name" : "Sentiment", "type" : "float", "null": true },
                { "name" : "ExtractedDates", "type" : "string", "null" : true },
                { "name" : "Date", "type" : "string", "default" : "" },
                { "name" : "CanDeleteWhenOld", "type" : "bool", "default" : true },

                { "name" : "Geo", "type" : "float_pair", "null" : true }
        ]
"joins" : [
{ "name" : "hasSource", "type" : "field", "store" : "Source", "inverse" : "hasSocialMediaInput" },
{ "name" : "hasConcept", "type" : "index", "store" : "Concept", "inverse" : "hasSocialMediaInput" },
{ "name" : "hasCategory", "type" : "index", "store" : "Category", "inverse" : "hasSocialMediaInput" },
{ "name" : "hasDate", "type" : "index", "store" : "Date", "inverse" : "hasSocialMediaInput" },
],
"keys" : [
                { "field" : "Title", "type" : "text", "vocabulary" : "text" },
                { "field" : "Body", "type" : "text", "vocabulary" : "text" },
                { "field" : "Language", "type" : "value" },
                { "field" : "Date", "type" : "value" }
        ],
"timeWindow" : {
        "duration" : 365,
        "unit" : "day",
        "field" : "DateTime"
        }
},
{
```

# User Interface

o **TWITTEROBSERVATORY** provides a suitable user interface that allows user to:

- **view** upcoming social media data (tweets),

- **search** tweets by different queries and

- **analyze** the search results within different dimensions.
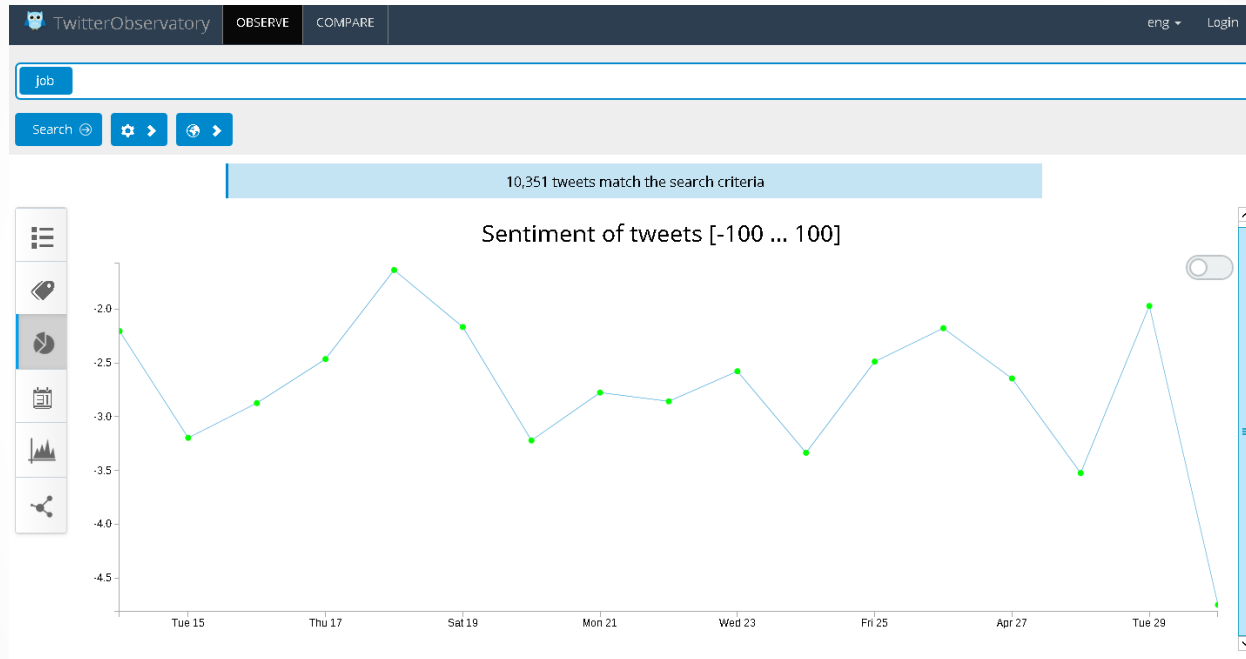
# User Interface: Observed Tweets



Observed Tweets with Details (Filter: "job")

# User Interface:
# Tag Cloud



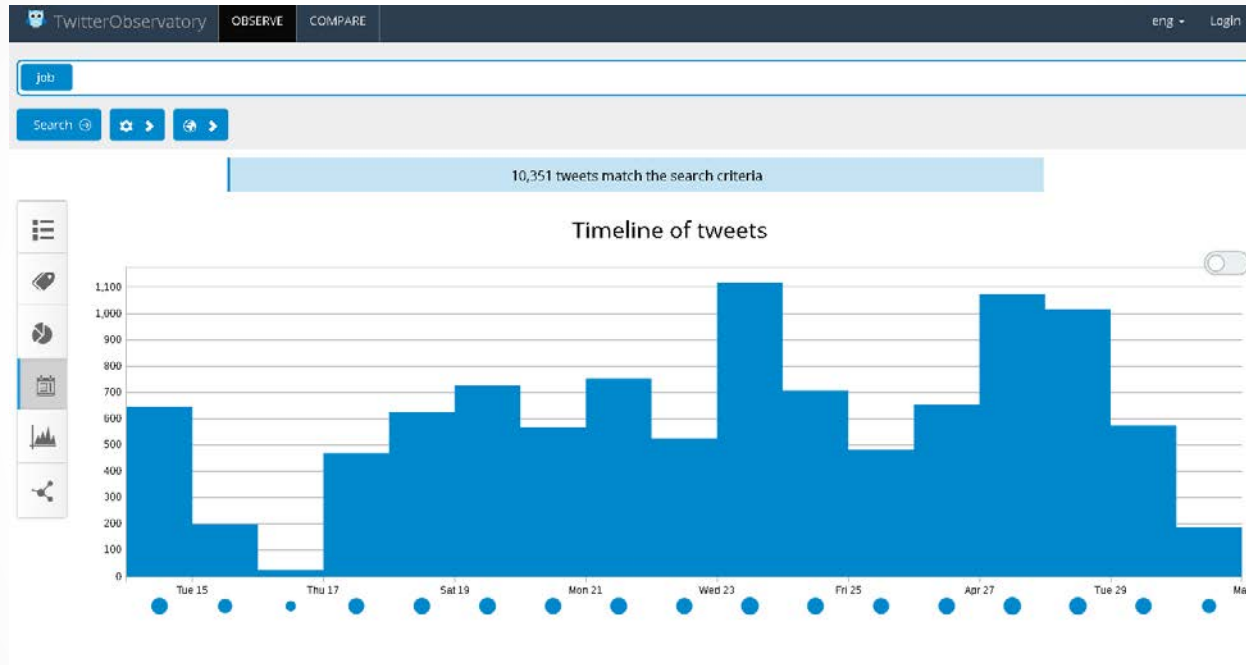Tag Cloud for Tweets (Filter: "job")

# User Interface: Sentiment



Sentiment for Tweets (Filter: "job")

# User Interface: Timeline



Tweets Timeline (Filter: "job")

# Introduction to Data Modeling

o **MODELING AND NOWCASTING** functionalities are intended to connect social media with external datasets, such as macroeconomic data.

o **Goal** of modeling and nowcasting is to **relate micro-signals** coming from social media (such as micro-signals related to stocks, micro-signals related to labor, micro-signals related to consumers, micro-signals related to real estate and credit, micro-signals related to energy) with **macro-economic variables**.

o **First test** - on data such as NTSF indices and other stock indices relevant to regional based crawling of tweets.

o **Combined features** from social media - **correlated** with macroeconomic time series, with a number of operators for time series analysis used (**moving average** (MA), **exponential moving average** (EMA), etc.

# Conclusion

o   We presented an **approach for social media mining** based on a pipeline that implements observing, enriching, storing, modeling and presentation techniques.

o   A novel tool **TwitterObservatory** that allows observing, searching, analyzing and presenting social media has been introduced.

o   The developed **software components** enable **monitoring of social media stream including enrichment and storing of the data**.

o   The **future work** will be based on implementing **additional functionalities for social media mining pipeline and on developing extensive modeling and nowcasting functionalities** for social media and external datasets.

# Questions?