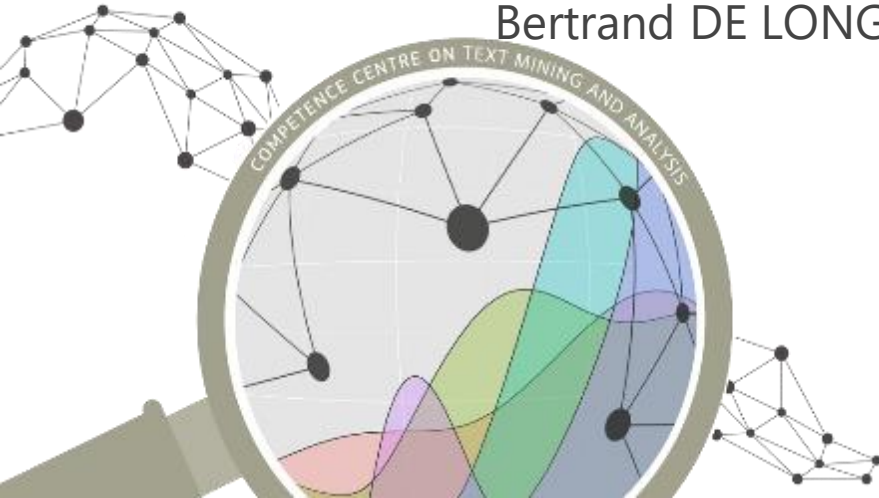# JRC's Text Mining and Analysis Competence Centre

Natural Language Processing for EU policy making

Bertrand DE LONGUEVILLE, DG JRC.I.3

The European Commission's science and knowledge service

Joint Research Centre

"To play a central role in creating, managing and making sense of collective scientific knowledge for better EU policies"

Petten

Geel

Karlsruhe

Brussels

Seville

Ispra

"Document Deluge"
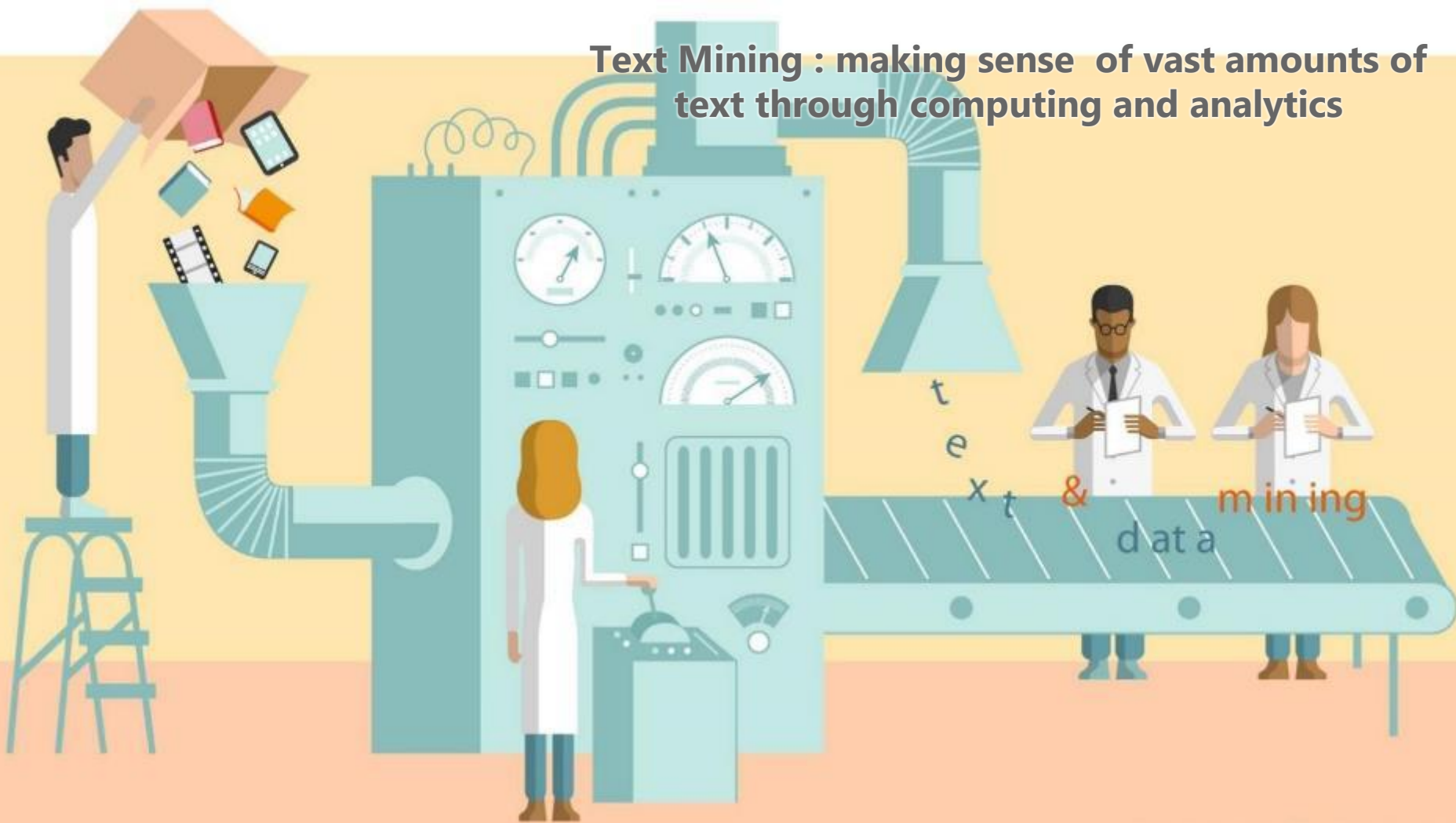
Climbing the pyramid

Wisdom

Knowledge

Information

Data

(CC0)

Text Mining : making sense of vast amounts of text through computing and analytics

*The Text Mining & Analysis Competence Centre is an in-house **consultancy** and **innovation** service supporting EU Institution's **Policy Makers, Investigators and Analysts** in their **knowledge-intensive tasks** by providing **advice** and advanced analytical **tools** in the field of text mining.*

# Some typical Text Mining use cases

# Sorting (aka "classifiers")

copyright © Rixie - Stock.Adobe.com

# Named Entity Recognition (NER)

In fact, the `Chinese` NORP market has the `three` CARDINAL most influential names of the retail and tech space – `Alibaba` GPE , `Baidu` ORG , and `Tencent` PERSON (collectively touted as `BAT` ORG ), and is betting big in the global `AI` GPE in retail industry space . The `three` CARDINAL giants which are claimed to have a cut-throat competition with the `U.S.` GPE (in terms of resources and capital) are positioning themselves to become the 'future `AI` PERSON platforms'. The trio is also expanding in other `Asian` NORP countries and investing heavily in the `U.S.` GPE based `AI` GPE startups to leverage the power of `AI` GPE .

Backed by such powerful initiatives and presence of these conglomerates, the market in APAC AI is forecast to be the fastest-growing `one` CARDINAL , with an anticipated `CAGR` PERSON of `45%` PERCENT over `2018 - 2024` DATE .

To further elaborate on the geograph... in `2017` DATE and has been leading, credit in the regional trends with ov... artificial intelligence technology. Add... such as `Google` ORG , `IBM` ORG , ...

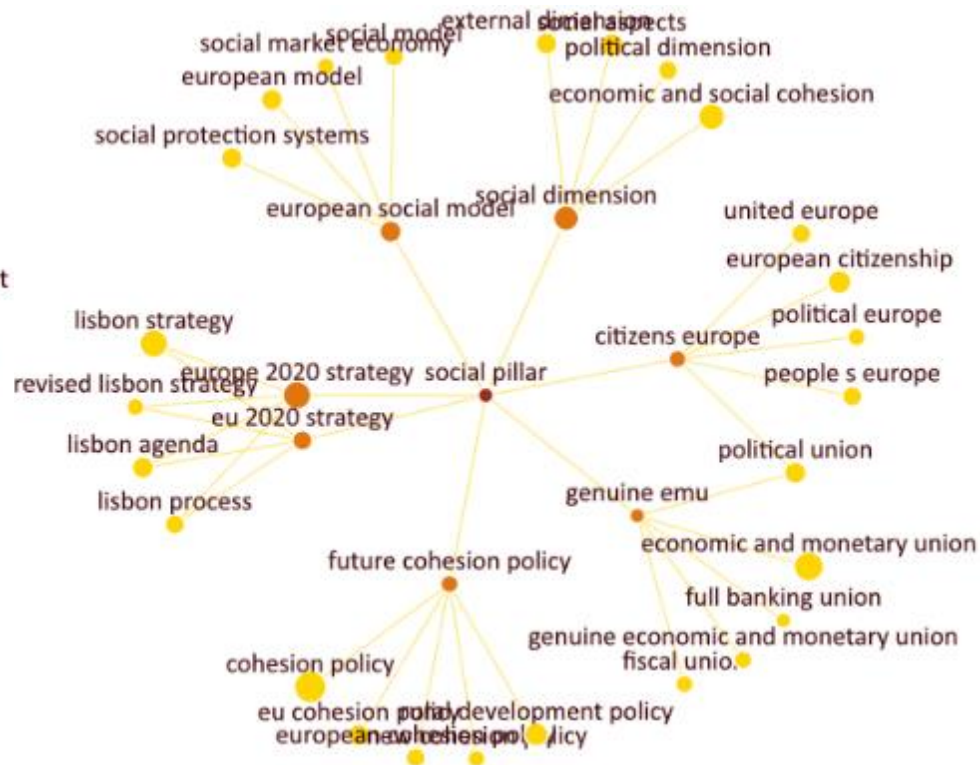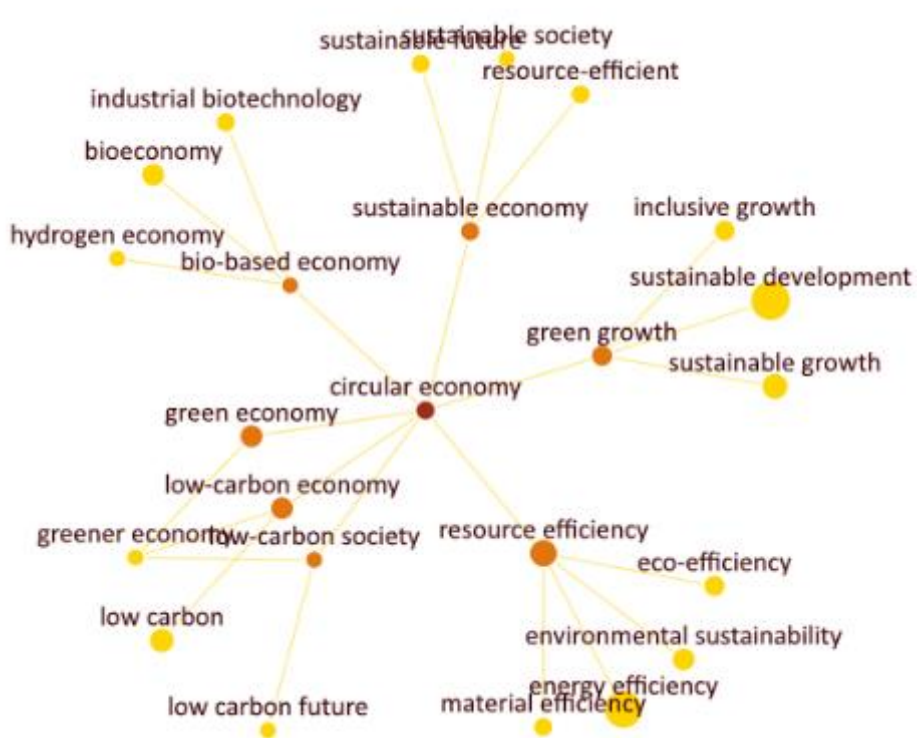| Arabic example | English translation | Entity type |
|---|---|---|
| أندونيسية / أندونيسب | Indonesia | Location |
| لوس انجيلس / لوس انجلوس / لوس انجليس / لوس انجيليس | Los Angeles | Location |
| جوهانسوبورغ/ جوهانسبرغ / جوهانسبورغة / جوهانسورغ | Johannesburg | Location |
| جلدر / جيلدر / غلدر / غيلدر | Guilder | Price (currency) |
| رقم الموبيل: ٥٧٥٦٤٥٣ ,الجوال: ٥٧٥٦٤٥٣ | Mobile no. 3546575 | Phone no. |

# Event Extraction

# Semantic Vector Space Analysis

# Tonality / Sentiment / Emotion detection

| Emotional Criteria | Example topic sentences |
| --- | --- |
| Trust | "Forbes Article Predicts Bitcoin Value will "Explode"" / "Good news for the Bitcoin" / "Don't panic, China is NOT banning bitcoin" |
| Fear | "Mining cartel attack" / "OMG! What has Satoshi created? He has opened Pandora's box" / "We are victims of our own success" |
| Surprise | "Whatever happened to the Bitcoin Police?" / "I think the rapture happened.....?" / "Blockchain.info "firstbits" changing/disappearing?!" |

# Flagship projects

A tool to monitor, analyse and aggregate online and other electronically available sources

## Wide Coverage

- Local, Regional, National and International
- Multi-Lingual

## Rapid Updates

- High Frequency (5min)

## Sources Covered

- Online media
- News wires
- Blogs
- BBC Monitoring
- Oxford Analytica
- Etc…

## Full Text Analysis

Europe Media Monitor

**Related projects**

**Many, many users worldwide**

MEDISYS

EMM Open Source Intelligence Suite

EU Institutions
Member States
International

# INNOVATION MONITORING

**~ 65 M documents**

# Data intelligence tools

$$\begin{bmatrix} Field_1 & Values_1 \\ \vdots & \vdots \\ Field_x & Values_x \end{bmatrix}$$

Data

Enrichment

Query

Index

Visualisation

**Visualisation**

**Pre-emerging Technologies detection**

# Semantic Analysis of Legal Texts

# The Power of Emotions on the (social) Web

-

**when data scientists aim for mindfulness in politics**

Matthew KING
Alexandra BALAHUR
Guillaume JACQUET
Bertrand DE LONGUEVILLE

Image credits : WikiMedia Commons

Image credits : WikiMedia Commons


Image credits : © Mounir via Adobe Stock

(social) Web Mining

+

Policy Making

=

?!

# emotions politics



Image credits : © European Parliament/Michel CHRISTEN

Image credits : JRC Report "Understanding our Political Nature"

Image credits : © AFP

Image credits : WikiMedia Commons

# warning : troll factory contents

# online vs real world ?!



EUSSR

**Percentage of EU citizens who tend to trust :**

**-The European Union**
**-Their National Government**
**-Their National Parliament**

Source : European Commission – Eurobarometer 91 – Spring 2019

EMOTION

# The Rhetorical Triangle

**Logos**
Logic. Reason. Proof.

**Pathos**
Emotions.
Values.

**Ethos**
Credibility.
Trust.

# what are we doing about this?

# From **Media Monitoring** to Public Opinion Analysis

Large Scale collection :
- 11,000 News Sites
- 325,000 articles per day
- 70 Languages
- 24/7 updates

Advanced inforation extraction
- Named Entity Recognition
- Quotes Extraction
- Events Extraction
- **Sentiment Analysis**

✅ **Available Right Now**

# From Media Monitoring to **Public Opinion Analysis**

**Top News Stories with Sentiment in end October 2019 ...**

EMM Large Scale News Corpus :
- 10+ years history
- About 1 Billion articles available for analysis

Media Analytics Capability:
- Sentiment per news story
- Sentiment per Person, organisation, topic, …
- Sentiment country/lanaguage
- **Combination of all of the above!**

**Work in Progress**

**... in Germany**

| | |
|---|---|
| | **Brexit** |
| | **Turkish offensive in Syria** |
| | **Halle (Germany), suspect arrested** |
| | **German political coalition** |
| | **Protests in Catalonia** |
| | **Merkel optimistic about deal** |

(Legend pie chart: neutral, joy, sad, disgust, anger, fear)

**... in Italy**

| | |
|---|---|
| | **Italian Budget** |
| | **Turkish offensive in Syria** |
| | **Brexit** |
| | **Whirlpool closing site** |
| | **Policemen Funerals in Trieste** |
| | **Hongkong riots** |

# From **Social Media Analytics** to Argument Mining



Social Media Analytics:
- Richer information on topics (MyNews)
- Fight Disinformation (MAT)
- Natural Disaster Management via rapid Tweets collection

✅ **Available Right Now**

# From Social Media Analysis to **Argument Mining**

The vision :
Fine-grained understanding of public debates, through strengthened JRC capabilities on Sentiment Analysis and Argument Mining

Way forward :
- Generalise successful 'one shot' experiments into a stable system
- Include more soucres of public debate : official communication of political actors, eurobarometer, public consultations (e.g. Future of Europe)

**Work in Progress**

**Volume of Tweets related to an EU policy and expressing a given sentiment, over time**



- anger
- sad
- joy
- surprise
- disgust
- fear

**Volume of Tweets arguing against EU migration policy, by argument type, over a given period**

Way forward : our research agenda

## Corpus Organisation

### A. Content categorisation and clustering

*Why it matters* enriching documents with topic information and organise news in cross-lingual stories are essential to most analysis use cases

*Where we are* EMM Category Matcher is in PROD, CLiCL on its way to PROD

*Where we go* redesign category concept for EMM 2.0, put CLiCL in prod, use it as an input for existing and new analyses

## Information Extraction

### B. Named Entities and Quotes Extraction

*Why it matters* extracting entities are essential to most analysis use cases; quotes are important for political analysis and Disinfo

*Where we are* NERONE, GEO and Quotes are in PROD (but not all languages equally covered)

*Where we go* improve recall, disambiguation and multilingualism, enhance NERONE for Knowledge Graph Extraction (see Research Area F.)

## Meaning Extraction

### C. Events extraction

*Why it matters* high demand for security (INTCEN, HOME, AU, *etc.*), epidemics (Medisys clients), disaster relief (ECHO); no solution from market

*Where we are* NEXUS is in PROD for 11 languages; quality is state-of-the-art (i.e. far from perfect)

*Where we go* experiment Machine Learning and Self-Learning (see Research Area H), more languages and event types, use CLiCL output

### D. Sentiment Analysis and Argument Mining

*Why it matters* high demand for policy analysis (SG, JRC, Policy DGs) and for Disinfo (COMM, LP, *etc.*)

*Where we are* SEmo in PROD for EN (FR, DE, IT, ES, PT are ready), prototypes developed for Argument Mining and Claim Checking

*Where we go* test SEmo with BERT, extend to EU-24 languages at least, resume Argument/Claims works based on Disinfo requirements

### E. Semantic Analysis of Texts

*Why it matters* enabler for topic mining, semantic search and Policy Analysis, Systematic Scientific Review, translation, summarisation and more

*Where we are* SeTA very close to PROD

*Where we go* build Policy Analysis prototype on top of SemLEX, test new "dialects" use cases based on input from OP, TIM (SemTECH), SG (SemPOL) and Social Media Mining activities (SemSOC)

## New Paradigms

### F. Knowledge Graph Extraction and Summarisation

*Why it matters* summarisation is requested by EMM clients (information overload); Knowledge Graph Extraction (KGE) is a cornerstone of EMM 2.0 vision (i.e. EMM as an Intelligence Analysis platform)

*Where we go* implement Extractive Summarisation as a quick win for EMM; lay foundations of KGE in EMM (to start with, enhance NERONE to this end)

### G. Multi-modal Text Analysis

*Why it matters* expected advances for Social Media Mining (i.e. analyse full context incl. media + social graph) and for Text Analysis combined with Big Data

*Where we go* define specific use case requirements, e.g. for Disinfo, Political Discourse Analysis, Emotional Landscape Analysis, etc.

### H. Self-leaning Models

*Why it matters* Machine Learning needs massive annotated data to increase performance; user feedback can provide such annotations at no cost

*Where we go* test this paradigm on systems for which feedback loops can be put in place, e.g. Event moderation, SEmo, NERONE (disambiguation)

### I. Prediction based on past news

*Why it matters* Foresight is a key strategic topic for JRC and EC Collège; EMM archive is a key asset for identifying weak signals announcing past events

*Where we go* get inspired by similar approaches (e.g. applied to financial stock market forecasts) ; identify a proper use case (e.g. election results forecast?) which can be tested with EMM archive

# Thank you for your attention

# questions

## are warmly welcome

European Commission - Joint Research Centre - Competences Directorate -Text and Data Mining Unit

bertrand.de-longueville@ec.europa.eu