# From data to journalism

Nelly Barret

4th year PhD student
Supervised by Ioana Manolescu
Inria Saclay and Institut Polytechnique de Paris

January 31, 2024

*Inria* · INSTITUT POLYTECHNIQUE DE PARIS · Le Monde · radiofrance · WEDODATA

# Personal (small) presentation

My background:

- CS Bachelor @ University of Lyon
- CS Master, AI track @ University of Lyon
- CS PhD student @ Inria and Ecole Polytechnique

My thesis is about **user-oriented exploration of semi-structured data**.

It is not only about me:



...and many others!

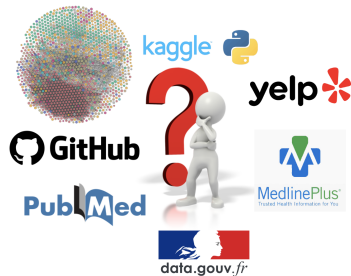Our digital world comes with various contexts, needs, actors, ...

We are **overwhelmed** by (raw) data, we need to bring order

Very **large** and **heterogeneous** data:

- Tables, text, databases, ...

Detection of **entities** of interest:

- People names, places, company names, dates, ...

# Data + journalism = data journalism
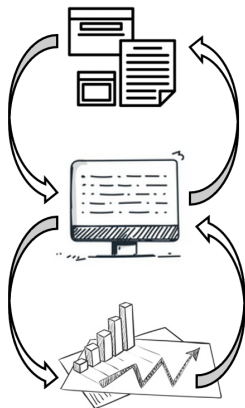
On one hand, we have:

- Facts
- Data

In the middle, we have:

- Computers
- Programs

On the other hand, we have:

- Journalists
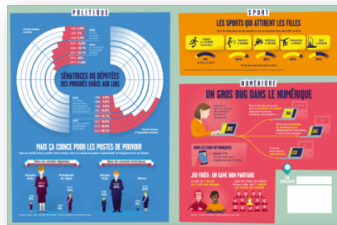- Data investigation
- Fact-checking

# New user needs, especially in data-journalism

## With **heterogeneous** data, users need:

1. A **uniform**/integrated view over the data
2. **Efficient** and **intuitive** ways to:
   - Get a global **understanding, description** of the data
   - Get interesting **entity connections**
   - **Query** and **search** for information in the system
3. Produce **insights**/**tangible results** to share

**But:**

- They have few or no CS skills
- They do not know what they are exactly looking for
- Their data may be messy/dirty

# Vocabulary introduction: dataset, schema, model

## Dataset
A <u>file</u> reporting data on a precise <u>topic</u>

## Data model
How data is <u>represented</u> (table, text, database, ...)

## Schema
How <u>data objects</u> are designed and relate

| produit | marque | genre | prix |
|---------|--------|-------|------|
| chemise | guess | homme | 50,99 |
| chaussure | adidas | femme | 44 |
| parfum | dior | femme | 120 |
| chemise | h&m | homme | 45 |

## Data heterogeneity
At the <u>model</u> and/or <u>schema</u> level

efficient and expressive integration of heterogeneous data

=

Provide a unified access (= put in the same "box")

efficient and expressive integration of heterogeneous data

=

Provide a unified access (= put in the same "box") to a set of datasets
(whatever their provenance, model, schema)

# Problem statement

efficient and expressive integration of heterogeneous data

=

Provide a unified access (= put in the same "box") to a set of datasets
(whatever their provenance, model, schema), sometimes very large

# Problem statement

efficient and expressive integration of heterogeneous data

=

Provide a unified access (= put in the same "box") to a set of datasets (whatever their provenance, model, schema), sometimes very large, such that this can be understood and used by human users.

# Problem statement

efficient and expressive integration of heterogeneous data
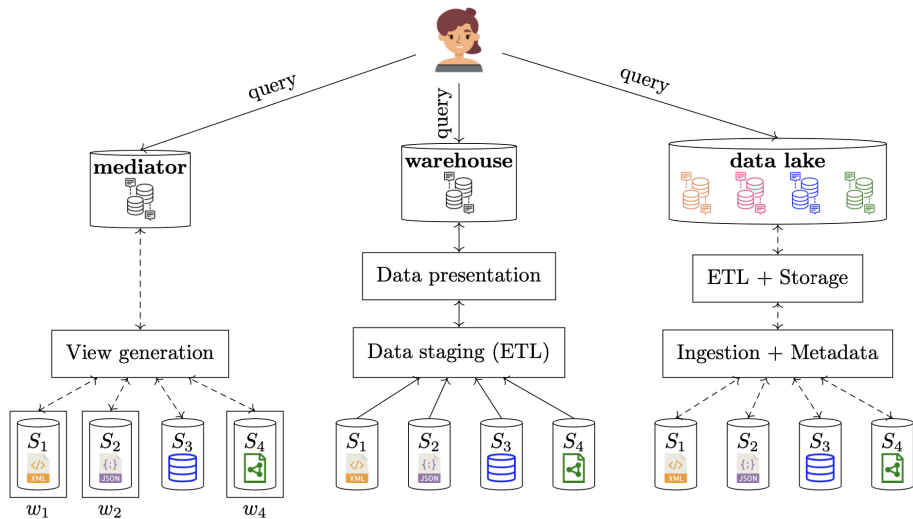
=

Provide a unified access (= put in the same "box") to a set of datasets (whatever their provenance, model, schema), sometimes very large, such that this can be understood and used by human users.

## What is data integration?

A system providing a unified interface to access, process and query a set of diverse, and potentially heterogeneous, datasets

# Existing architectures for data integration

## Yes, but...

Data integration systems strengths are also their weaknesses:

- **Mediators** convert many sources to a single model, but...
    - $\rightarrow$ Not feasible for dozens of sources
- **Warehouses** lead to a consolidated database, but...
    - $\rightarrow$ Not very flexible with new data
- **Data lakes** allow many data sources to co-exist, but...
    - $\rightarrow$ Rapidly become data swamps

No data integration systems fits all needs!
Data integration takes time, money and requires CS skills

# Our proposals: ConnectionStudio and StatCheck

## ① ConnectionStudio

- A <u>data lake</u> for <u>novice users</u>
- To load, clean, visualize and query heterogeneous data
- → With "LeMonde" data journalists
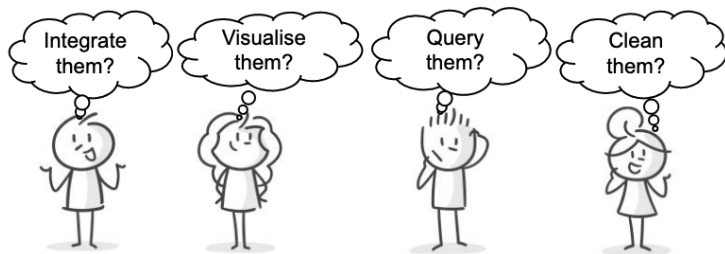
## ② StatCheck

- A <u>warehouse</u> for centralizing <u>statistical data</u>
- To search for statistics and to analyse political discourses
- → With the "FranceInfo" fact-checking team "Le vrai du faux"
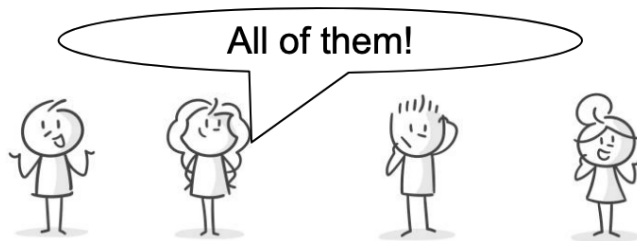
# ConnectionStudio

# ConnectionStudio problem statement

**Question:** How to work with heterogeneous (journalistic) data?

**Question:** How to work with heterogeneous (journalistic) data?

# ConnectionStudio problem statement

**Question:** How to work with heterogeneous (journalistic) data?

# ConnectionStudio problem statement

**Question:** How to work with heterogeneous (journalistic) data?



All of them!

# ConnectionStudio solution in a nutshell

**Our answer:** A **user-friendly** <u>data lake</u> for:

# ConnectionStudio solution in a nutshell

**Our answer:** A **user-friendly** <u>data lake</u> for:

① **Loading heterogeneous datasets**
→ uniform/integrated view over the data

# ConnectionStudio solution in a nutshell

**Our answer:** A **user-friendly** <u>data lake</u> for:

① **Loading heterogeneous datasets**
  → uniform/integrated view over the data
② **Statistics and data summaries**
  → global description of the data

# ConnectionStudio solution in a nutshell

**Our answer:** A **user-friendly** <u>data lake</u> for:

① **Loading heterogeneous datasets**
   → uniform/integrated view over the data

② **Statistics and data summaries**
   → global description of the data

③ **Keyword search**
   → query and search for information

# ConnectionStudio solution in a nutshell

**Our answer:** A **user-friendly** <u>data lake</u> for:

① **Loading heterogeneous datasets**
  $\rightarrow$ uniform/integrated view over the data

② **Statistics and data summaries**
  $\rightarrow$ global description of the data

③ **Keyword search**
  $\rightarrow$ query and search for information

④ **Data paths**
  $\rightarrow$ interesting entity connections
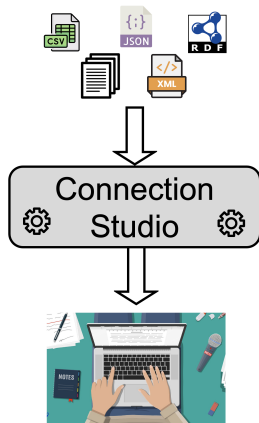
# ConnectionStudio solution in a nutshell

**Our answer:** A **user-friendly** <u>**data lake**</u> for:

① **Loading heterogeneous datasets**
   → uniform/integrated view over the data

② **Statistics and data summaries**
   → global description of the data

③ **Keyword search**
   → query and search for information

④ **Data paths**
   → interesting entity connections

⑤ **Querying the data lake**
   → query and search for information

# ConnectionStudio solution in a nutshell

**Our answer:** A **user-friendly** <u>data lake</u> for:

① **Loading heterogeneous datasets**
  → uniform/integrated view over the data

② **Statistics and data summaries**
  → global description of the data

③ **Keyword search**
  → query and search for information

④ **Data paths**
  → interesting entity connections

⑤ **Querying the data lake**
  → query and search for information

⑥ **Tabular-looking results**
  → produce insights/tangible results to share



Connection Studio

**Question:** How to commonly represent heterogeneous datasets?

# ① Unified data view: a graph

**Question:** How to commonly represent heterogeneous datasets?

**Our answer:** A graph

# ① Unified data view: a graph

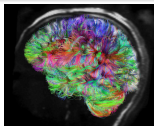**Question:** How to commonly represent heterogeneous datasets?

**Our answer:** A graph

## Wait... What is a graph?

The **graph** paradigm describes:

- Objects (nodes)
- Connected by links (edges)
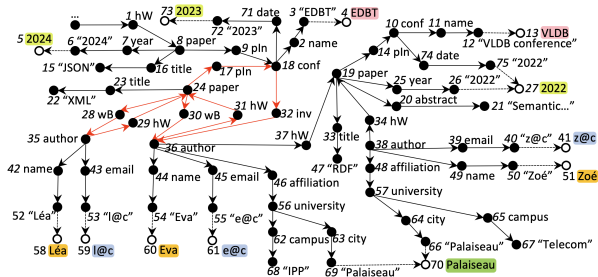
High flexibility → largely used



Brain neurons



International flights



Panama papers

# ① Unified data view: a graph

- Ingest any dataset into a **directed graph** (•, →)
- Extract **named entities**, NEs, from the graph values (○, --→):
  - Temporal: date , time reference
  - Web: URI, email address , hashtag, Twitter citation
  - Complex entities: People   Place   Organization
  - Used pre-trained language models; more recently ChatGPT

**Question:** What if the data graph is huge?

$\left(1\right)$ Unified data view: a graph
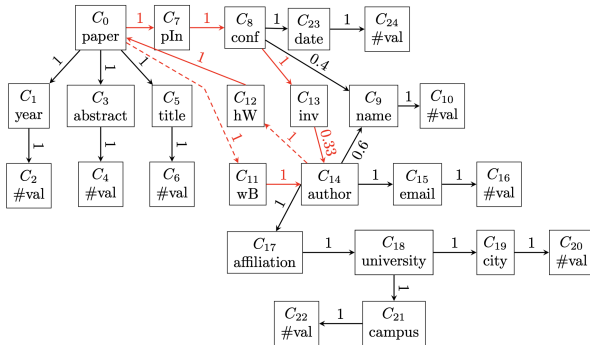
**Question:** What if the data graph is huge?

**Our answer:** Build its compact representation (summary)

**Question:** What if the data graph is huge?

**Our answer:** Build its compact representation (summary)

→ We build a **summary graph**, with small information loss
→ **Efficient** algorithms and applications

**Question:** How to get a quick overview of the datasets?

# ② Statistics and data summaries

**Question:** How to get a quick overview of the datasets?

**Our answer:** Show Named Entities stats in charts, tables and tag clouds

**Question:** How to get a quick overview of the datasets?

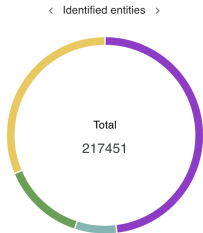**Our answer:** Show Named Entities stats in charts, tables and tag clouds

**Question:** How to get a quick overview of the datasets?

**Question:** How to get a quick overview of the datasets?

**Our answer:** Draw an Entity-Relationship diagram representing "important" nodes

**Question:** How to get a quick overview of the datasets?

**Our answer:** Draw an Entity-Relationship diagram representing "important" nodes

Abstraction of conferences
located at: file:/Users/nelly/Documents/boulot/theseNelly/abstraction-work/../abstraction-data/rdf/testConferences.nt
with 103 normalized nodes, 29 collections, (PR, FLAC), 2/5 main collections, data coverage is 100.0%



4 Author (http://dbpedia.org/ontology/writer) ⓘ
- affiliation
  - University
    - name
- orcid
- teaching
- field
- name

2 Conference (Event)
- name
- duration
- place

attends
publishes.Paper.writtenBy

Abstraction of xmark00125
located at: file:/Users/nelly/Documents/boulot/theseNelly/abstraction-work/../abstraction-data/xml/xmark0.0125.xml
with 42293 normalized nodes, 136 collections, (PR, FLAC), 5/5 main collections, data coverage is 91.0%

**Question:** How to "discuss"/ask something to the data lake?

**Question:** How to "discuss"/ask something to the data lake?

Language to "discuss"/ask something to a database: (mostly) SQL

**Question:** How to "discuss"/ask something to the data lake?

Language to "discuss"/ask something to a database: (mostly) SQL

```
SELECT n1.label, n2.label, n3.label
FROM nodes n1, edges e1,
     nodes n2, edges e2, nodes n3
WHERE e1.source=n1.id
      AND e1.target=n2.id
      AND e2.source=n2.id
      AND e2.target=n3.id
LIMIT 3
```

| n1.label | n2.label | n3.label |
|----------|----------|----------|
| author | name | "Léa" |
| university | campus | "IPP" |
| paper | writtenBy | author |

**Question:** How to "discuss"/ask something to the data lake?

Language to "discuss"/ask something to a database: (mostly) SQL

```
SELECT n1.label, n2.label, n3.label
FROM nodes n1, edges e1,
     nodes n2, edges e2, nodes n3
WHERE e1.source=n1.id
      AND e1.target=n2.id
      AND e2.source=n2.id
      AND e2.target=n3.id
LIMIT 3
```

| n1.label | n2.label | n3.label |
|----------|----------|----------|
| author | name | "Léa" |
| university | campus | "IPP" |
| paper | writtenBy | author |

Pros and cons of SQL:

- Needs to be learned ($\rightarrow$ time, skills)
- SQL queries are highly performant, scalable, optimizable
- Results shown as tables

**Our answer:** Get rid of the "SQL writing part", keep tables as output

**Our answer:** Get rid of the "SQL writing part", keep tables as output

| Path 1 | | Starting variable | Ending variable | |
|---|---|---|---|---|
| declaration.general.declarer.name#val | | decla | deputyName | ● EVALUATE THE QUERY   ⬆ SAVE CHANGES |

| Path 2 | | Starting variable | Ending variable | Join |
|---|---|---|---|---|
| declaration.financialInterest.items.item | | decla | item | ● Required ○ Optional 🗑 |

| Path 3 | | Starting variable | Ending variable | Join |
|---|---|---|---|---|
| item.company#val.extract:o | | item | companyName | ● Required ○ Optional 🗑 |

| Path 4 | | Starting variable | Ending variable | Join |
|---|---|---|---|---|
| item.nbShares#val | | item | nbShares | ○ Required ● Optional 🗑 |

| Path 5 | | Starting variable | Ending variable | Join |
|---|---|---|---|---|
| row.company_name.#val.extract:o | | csvline | companyName | ● Required ○ Optional 🗑 |

**Ⅲ COLUMNS   ☰ FILTERS   ☰ DENSITY   ⬆ EXPORT**

| decla | deputyname | item | companyname | nbshares | csvline |
|---|---|---|---|---|---|
| 2660 | alain pierre marie rousset | 2743 | sanofi | 1200 | 352 |
| 1470 | edouard courtial | 1511 | lvmh | 29013 | 248 |
| 1470 | edouard courtial | 1543 | michelin | 162179 | 261 |

**Our answer:** Get rid of the "SQL writing part", keep tables as output

**Our answer:** Get rid of the "SQL writing part", keep tables as output

1. Enumerate a set of paths in the summary graph:
   - $p = \{n_0, e_0, n_1, e_1, ..., e_i, n_{i+1}\}$
2. Each selected path $p$ is associated to:
   - A source variable $s$ (the first element in $p$)
   - A target variable $t$ (the last element in $p$)
3. Select join predicates (`LEFT JOIN` or `INNER JOIN`)
4. Conversion to a SQL query:
   - Each path leads to a SQL query, reusing $s$ and $t$
   - Path SQL queries are joined using join predicates

$$p_0 = \{\overbrace{declaration}^{s}, \_, general, \_, declarer, \_, name, \_, \overbrace{\#val}^{t}\}$$
$$\bowtie$$
$$p_1 = \{\underbrace{declaration}_{s}, \_, financialInterest, \_, items, \_, \underbrace{item}_{t}\}$$

**Question:** How to share results found/created in the data lake?

## ⑥ Concrete results

**Question:** How to share results found/created in the data lake?

**Our answer:** They can be exported (statistics, tables, queries, ...)

**Question:** How to share results found/created in the data lake?

**Our answer:** They can be exported (statistics, tables, queries, ...)

# StatCheck

# StatCheck problem statement

**Question:** How to ease/automate part of the fact-checking process?

# StatCheck problem statement

**Question:** How to ease/automate part of the fact-checking process?

# ConnectionStudio problem statement

**Question:** How to ease/automate part of the fact-checking process?

# ConnectionStudio problem statement

**Question:** How to ease/automate part of the fact-checking process?

## StatCheck in a nutshell

**Our answer:** A **user-friendly** <u>warehouse</u> for:

# StatCheck in a nutshell

**Our answer:** A **user-friendly** <u>warehouse</u> for:

①  Consolidating French statistical data
    $\rightarrow$ a unique repository for French statistics

# StatCheck in a nutshell

**Our answer:** A **user-friendly** <u>**warehouse**</u> for:

(1) Consolidating French statistical data
   → a unique repository for French statistics

(2) Statistical search engine
   → search for statistical information

# StatCheck in a nutshell

**Our answer:** A **user-friendly** <u>warehouse</u> for:

① Consolidating French statistical data
   $\rightarrow$ a unique repository for French statistics

② Statistical search engine
   $\rightarrow$ search for statistical information

③ Automated text analysis 1/2
   $\rightarrow$ detection of statistical claims

# StatCheck in a nutshell

**Our answer:** A **user-friendly** <u>warehouse</u> for:

① Consolidating French statistical data
  → a unique repository for French statistics

② Statistical search engine
  → search for statistical information

③ Automated text analysis 1/2
  → detection of statistical claims

④ Automated text analysis 2/2
  → recognise persuasion techniques in political discourses

**Question:** how to create consolidated data for French statistics?

**Question:** how to create consolidated data for French statistics?

**Our answer:** crawling of reference/trusted websites for French statistics (Open Data):

# 1 Consolidated data for French statistics

**Question:** how to create consolidated data for French statistics?

**Our answer:** crawling of reference/trusted websites for French statistics (Open Data):

1. **INSEE**: Institut National de la Statistique et des Études Économiques
2. **EuroStat**: European statistical database

**Question:** how to create consolidated data for French statistics?

**Our answer:** crawling of reference/trusted websites for French statistics (Open Data):

1. **INSEE**: Institut National de la Statistique et des Études Économiques
2. **EuroStat**: European statistical database

Knowing that:

- Data is in different models (tables, text, ...)
- Data is huge
- ... and many other concerns that we will not cover today

## ① Consolidated data for French statistics

**Initial approach:** convert all data into a graph

**Yes, but:**

- High cost of storage (graph size: 3Tb)
- Searching the graph was expensive (1/3 of queries were very long)

**Initial approach:** convert all data into a graph

**Yes, but:**

- High cost of storage (graph size: 3Tb)
- Searching the graph was expensive (1/3 of queries were very long)

| Code région | Code EPCI | Code de l'unité urbaine | Libellé géographique | Nombre de logements du Parc Locatif Social | Nombre de logements sociaux mis en service dans l'année | Taux de vacance des logements sociaux | Taux de vacance de plus de 3 mois des logements sociaux | Taux de rotation des logements sociaux | Part des logements sociaux collectifs | Part des logements sociaux individuels | Part des logements sociaux d'une pièce | Part des logements sociaux de deux pièces |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Reg | EPCI | UU2020 | LibGeo | nbLsPls | nbLsMes | txVac | txVac3m | txRot | txLsCol | txLsInd | txLs1p | txLs2p |
| 84 | 240100883 | 1303 | Les Pérouses-Triangle d'Activités | 315 | | | | 17.3 | | | 18.7 | |
| 84 | 240100883 | 1303 | Longeray-Gare | 878 | | | | 9.6 | 94.6 | 5.4 | 5.0 | 14.5 |
| 84 | 240100883 | 1303 | Centre-Saint-Germain-Vareilles | 528 | | 4.4 | | 13.9 | | | | 18.2 |
| 84 | 240100883 | 1303 | Tiret-Les Allymes | 290 | | | | 7.1 | | | | |
| 84 | 240100891 | 360 | Centre Ville | | | | | | | | | |
| 84 | 240100891 | 360 | Lancrans-Coupy-Vanchy | | | | | | | | | |
| 84 | 240100891 | 360 | Arc Vouvray-Gare-Châtillon | | | | | | | | | |
| 84 | 240100891 | 360 | Plateau de Musinens | | | | | | | | | |
| 84 | 240100891 | 360 | Arlod | | | | | | | | | |
| 84 | 240100891 | 360 | Châtillon-en-Michaille | | | | | | | | | |
| 84 | 200040350 | 1301 | Ouest | 135 | | | | 13.3 | | | | 17.8 |
| 84 | 200040350 | 1301 | Centre et Est | 489 | 25 | 7.7 | 4.8 | 12.8 | | | | 22.5 |
| 84 | 200040350 | 1301 | Sud-Ouest | 507 | | | | 7.5 | 87.4 | 12.6 | | 18.3 |
| 84 | 200071751 | 1501 | Centre Ville | 137 | | | | | | | 25.5 | 27.0 |
| 84 | 200071751 | 1501 | Champ de Foire | 248 | | | | 7.6 | | | | 49.6 |
| 84 | 200071751 | 1501 | Préfecture | 234 | | 5.2 | | 19.3 | | | 12.0 | 23.1 |
| 84 | 200071751 | 1501 | Citadelle | 392 | | 4.3 | | 7.9 | 74.2 | 25.8 | | 28.8 |
| 84 | 200071751 | 1501 | Mail | 571 | | | | 13.2 | 93.0 | 7.0 | 5.4 | 18.2 |
| 84 | 200071751 | 1501 | Peloux | 295 | | 5.2 | 4.5 | 10.5 | 92.9 | 7.1 | 5.1 | 16.6 |
| 84 | 200071751 | 1501 | Gare | 221 | | | | 11.7 | | | 27.6 | |
| 84 | 200071751 | 1501 | Brou | 643 | | 2.4 | | 11.3 | | | 5.4 | 24.6 |

**Our answer:** convert data tables as "areas"

**Our answer:** convert data tables as "areas"

| Code région | Code EPCI | Code de l'unité urbaine | Libellé géographique | Nombre de logements du Parc Locatif Social | Nombre de logements sociaux mis en service dans l'année | Taux de vacance des logements sociaux | Taux de vacance de plus de 3 mois des logements sociaux | Taux de rotation des logements sociaux | Part des logements sociaux collectifs | Part des logements sociaux individuels | Part des logements sociaux d'une pièce | Part des logements sociaux de deux pièces |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Reg | EPCI | UU2020 | LibGeo | nbLsPls | nbLsMes | txVac | txVac3m | txRot | txLsCol | txLsInd | txLs1p | txLs2p |
| 84 | 240100883 | 1303 | Les Pérouses-Triangle d'Activités | 315 | | | | 17.3 | | | 18.7 | |
| 84 | 240100883 | 1303 | Longeray-Gare | 878 | | | | 9.6 | 94.6 | 5.4 | 5.0 | 14.5 |
| 84 | 240100883 | 1303 | Centre-Saint-Germain-Vareilles | 528 | | 4.4 | | 13.9 | | | | 18.2 |
| 84 | 240100883 | 1303 | Tiret-Les Allymes | 290 | | | | 7.1 | | | | |
| 84 | 240100891 | 360 | Centre Ville | | | | | | | | | |
| 84 | 240100891 | 360 | Lancrans-Coupy-Vanchy | | | | | | | | | |
| 84 | 240100891 | 360 | Arc Vouvray-Gare-Châtillon | | | | | | | | | |
| 84 | 240100891 | 360 | Plateau de Musinens | | | | | | | | | |
| 84 | 240100891 | 360 | Arlod | | | | | | | | | |
| 84 | 240100891 | 360 | Châtillon-en-Michaille | | | | | | | | | |
| 84 | 200040350 | 1301 | Ouest | 135 | | | | 13.3 | | | | 17.8 |
| 84 | 200040350 | 1301 | Centre et Est | 489 | 25 | 7.7 | 4.8 | 12.8 | | | | 22.5 |
| 84 | 200040350 | 1301 | Sud-Ouest | 507 | | | | 7.5 | 87.4 | 12.6 | | 18.3 |
| 84 | 200071751 | 1501 | Centre Ville | 137 | | | | | | | 25.5 | 27.0 |
| 84 | 200071751 | 1501 | Champ de Foire | 248 | | | | 7.6 | | | | 49.6 |
| 84 | 200071751 | 1501 | Préfecture | 234 | | 5.2 | | 19.3 | | | 12.0 | 23.1 |
| 84 | 200071751 | 1501 | Citadelle | 392 | | 4.3 | | 7.9 | 74.2 | 25.8 | | 28.8 |
| 84 | 200071751 | 1501 | Mail | 571 | | | | 13.2 | 93.0 | 7.0 | 5.4 | 18.2 |
| 84 | 200071751 | 1501 | Peloux | 295 | | 5.2 | 4.5 | 10.5 | 92.9 | 7.1 | 5.1 | 16.6 |
| 84 | 200071751 | 1501 | Gare | 221 | | | | 11.7 | | | 27.6 | |
| 84 | 200071751 | 1501 | Brou | 643 | | 2.4 | | 11.3 | | | 5.4 | 24.6 |

**Our answer:** convert data tables as "areas"

**Our answer:** convert data tables as "areas"

| Code région | Code EPCI | Code de l'unité urbaine | Libellé géographique | Nombre de logements du Parc Locatif Social | Nombre de logements sociaux mis en service dans l'année | Taux de vacance des logements sociaux | Taux de vacance de plus de 3 mois des logements sociaux | Taux de rotation des logements sociaux | Part des logements sociaux collectifs | Part des logements sociaux individuels | Part des logements sociaux d'une pièce | Part des logements sociaux de deux pièces |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Reg** | **EPCI** | **UU2020** | **LibGeo** | **nbLsPls** | **nbLsMes** | **txVac** | **txVac3m** | **txRot** | **txLsCol** | **txLsInd** | **txLs1p** | **txLs2p** |
| 84 | 240100883 | 1303 | Les Pérouses-Triangle d'Activités | 315 | | | | 17.3 | | | 18.7 | |
| 84 | 240100883 | 1303 | Longeray-Gare | 878 | | | | 9.6 | 94.6 | 5.4 | 5.0 | 14.5 |
| 84 | 240100883 | 1303 | Centre-Saint-Germain-Vareilles | 528 | | 4.4 | | 13.9 | | | | 18.2 |
| 84 | 240100883 | 1303 | Tiret-Les Allymes | 290 | | | | 7.1 | | | | |
| 84 | 240100891 | 360 | Centre Ville | | | | | | | | | |
| 84 | 240100891 | 360 | Lancrans-Coupy-Vanchy | | | | | | | | | |
| 84 | 240100891 | 360 | Arc Vouvray-Gare-Châtillon | | | | | | | | | |
| 84 | 240100891 | 360 | Plateau de Musinens | | | | | | | | | |
| 84 | 240100891 | 360 | Arlod | | | | | | | | | |
| 84 | 240100891 | 360 | Châtillon-en-Michaille | | | | | | | | | |
| 84 | 200040350 | 1301 | Ouest | 135 | | | | 13.3 | | | | 17.8 |
| 84 | 200040350 | 1301 | Centre et Est | 489 | 25 | 7.7 | 4.8 | 12.8 | | | | 22.5 |
| 84 | 200040350 | 1301 | Sud-Ouest | 507 | | | | 7.5 | 87.4 | 12.6 | | 18.3 |
| 84 | 200071751 | 1501 | Centre Ville | 137 | | | | 27.0 | | | 25.5 | 27.0 |
| 84 | 200071751 | 1501 | Champ de Foire | 248 | | | | 7.6 | | | | 49.6 |
| 84 | 200071751 | 1501 | Préfecture | 234 | | 5.2 | | 19.3 | | | 12.0 | 23.1 |
| 84 | 200071751 | 1501 | Citadelle | 392 | | 4.3 | | 7.9 | 74.2 | 25.8 | | 28.8 |
| 84 | 200071751 | 1501 | Mail | 571 | | | | 13.2 | 93.0 | 7.0 | 5.4 | 18.2 |
| 84 | 200071751 | 1501 | Peloux | 295 | | 5.2 | 4.5 | 10.5 | 92.9 | 7.1 | 5.1 | 16.6 |
| 84 | 200071751 | 1501 | Gare | 221 | | | | 11.7 | | | 27.6 | |
| 84 | 200071751 | 1501 | Brou | 643 | | 2.4 | | 11.3 | | | 5.4 | 24.6 |

**Question:** How much more efficient is the novel approach?

**Question:** How much more efficient is the novel approach?

**Our answer:**

# ① Consolidated data for French statistics

**Question:** How much more efficient is the novel approach?

**Our answer:**

| Type | INSEE | EuroStat | Total |
|------|------:|---------:|------:|
| Files | 96 207 | 7 094 | 103 301 |
| Tables | 112 966 | 7 003 | 119 969 |
| Areas | 1 286 603 | 12 179 533 | **13 488 136** |
| Graph (Mb) | 1 864 766 | 120 425 | - |
| Areas (Mb) | 577 | 8 055 | - |
| **Compression** | × **3 266** | × **14** | - |

# ① Consolidated data for French statistics

**Question:** How much more efficient is the novel approach?

**Our answer:**

| Type | INSEE | EuroStat | Total |
|------|------:|---------:|------:|
| Files | 96 207 | 7 094 | 103 301 |
| Tables | 112 966 | 7 003 | 119 969 |
| Areas | 1 286 603 | 12 179 533 | **13 488 136** |
| Graph (Mb) | 1 864 766 | 120 425 | - |
| Areas (Mb) | 577 | 8 055 | - |
| **Compression** | × **3 266** | × **14** | - |

**Lesson learned:** the storage should be chosen based on the data/usage

**Question:** How to retrieve information from textual questions?

**Question:** How to retrieve information from textual questions?

**Our answer:** Given a query $Q$ composed of $n$ keywords $k_n$:

**Question:** How to retrieve information from textual questions?

**Our answer:** Given a query $Q$ composed of $n$ keywords $k_n$:

$$Q = \underbrace{\text{Taux}}_{k_1} \text{ de } \underbrace{\text{chômage}}_{k_2} \underbrace{2022}_{k_3}$$

**Question:** How to retrieve information from textual questions?

**Our answer:** Given a query $Q$ composed of $n$ keywords $k_n$:

$$Q = \underbrace{\text{Taux}}_{k_1} \text{ de } \underbrace{\text{chômage}}_{k_2} \underbrace{2022}_{k_3}$$

Find the 20 most relevant tables, and possibly the value:



Taux de chômage au 1er trimestre 2022 par région métropolitaine (en %) en %
Publiée le 08 juillet 2022



Taux de chômage au 1er trimestre 2022 dans les départements normands
Publiée le 08 juillet 2022

**Question:** Can we learn things from political discourses?

**Question:** Can we learn things from political discourses?

**Our answer:** Yes, e.g., in tweets

**Question:** Can we learn things from political discourses?

**Our answer:** Yes, e.g., in tweets

- Recognize/extract <u>statistical claims</u>

- Identify well-known <u>persuasion techniques</u>

| | |
|---|---:|
| Number of followed politicians | 63 |
| Number of gathered tweets | 77 081 |
| Claims detected (since 01/2022) | 61 207 |

**Question:** How to identify statistics used in political discourses?

**Question:** How to identify statistics used in political discourses?

**Our answer:** Use a Machine Learning model trained on political debates (from 1950 to 2024)
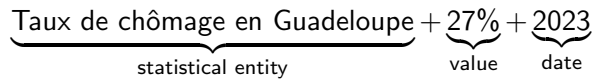
**Question:** How to identify statistics used in political discourses?

**Our answer:** Use a Machine Learning model trained on political debates (from 1950 to 2024)

- statistical claim = a statistical entity + a value + a date

$$\underbrace{\text{Taux de chômage en Guadeloupe}}_{\text{statistical entity}} + \underbrace{27\%}_{\text{value}} + \underbrace{2023}_{\text{date}}$$

**Marine Le Pen**, le 20 janvier 2023 à 16:47                    ⟨  •• ⟩  100%

Selon l' `ORG Insee` , `NUM 27%` des `CONT_ENT jeunes` en `LOC Guadeloupe` sont sans `ENT emploi` ni `CONT_ENT formation` . Il

est urgent de remettre l'Outre-mer au cœur des priorités et des politiques publiques, et de

créer les conditions qui favorisent l'investissement, gage de développement économique.

**Question:** Can we detect persuasion techniques in political discourses?

# ③ Persuasion techniques detection in political discourses

**Question:** Can we detect persuasion techniques in political discourses?

**Our answer:** Use ML classifiers to detect those techniques and assign them a category

**Question:** Can we detect persuasion techniques in political discourses?

**Our answer:** Use ML classifiers to detect those techniques and assign them a category

Emmanuel Macron on Nov 28, 2023:

Le cap que je porte a toujours été le même:

*affirmation*
réindustrialiser la France, gagner la bataille du plein-emploi,

*flag-waving*            *war term*
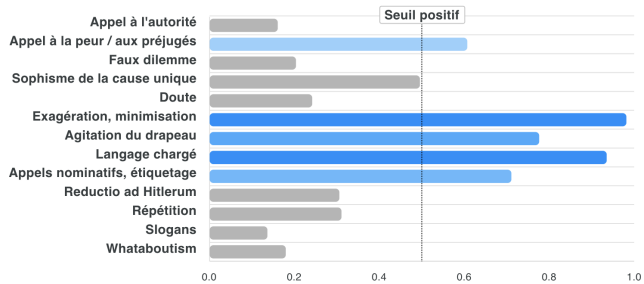être une Nation plus souveraine, industrielle et décarbonée.

*flag-waving, hopes*
Soyons des optimistes déterminés.

*loaded language*

1. Binary classifier:
   - Given a sentence, does it contain persuasive techniques?
2. Multi-class classifier:
   - Given a sentence, which of the 13 persuasive techniques are used?



Score du modèle de détection: 81.80%

# Takeaways and future work

**Takeaways:**

- **ConnectionStudio**: a user-oriented data lake for data exploration
- **StatCheck**: a statistical warehouse for fact-checking

**Future work:**

- ConnectionStudio:
    - Link data graph Named Entities to Wikipedia/trusted resources
    - Propose new ways to query the data
    - Clean (automatically) data in the data lake
- StatCheck:
    - Gathering other sources than INSEE and EuroStat
    - Cross-check statistical data between sources
    - Investigate recent Machine Learning models

# Final words

**ConnectionStudio**  **StatCheck**

If you are interested in what we are doing in the CEDAR team at Inria

# Next: interactive sessions

- 2 groups
- 1 hour each



We will put our journalists hats:

1. Investigate a use-case in ConnectionStudio
2. Browse StatCheck data, tweets and ML outputs

And discuss about your questions, thoughts, ...

# Next: interactive sessions

**CAC40**: CSV dataset
- Describes the top-40 most influential French companies
- Quite small (40 lines, 3 columns)

**HATVP**: XML dataset
- Describes political members' declarations about their wealth, jobs, financial intereset, ...
- Large ($\sim$ 2M nodes, $\sim$ 2M edges)

**They share Named Entities:**
- Crédit agricole, Danone, Education Nationale, Bouygues, ...

Let's see what we can do in ConnectionStudio!