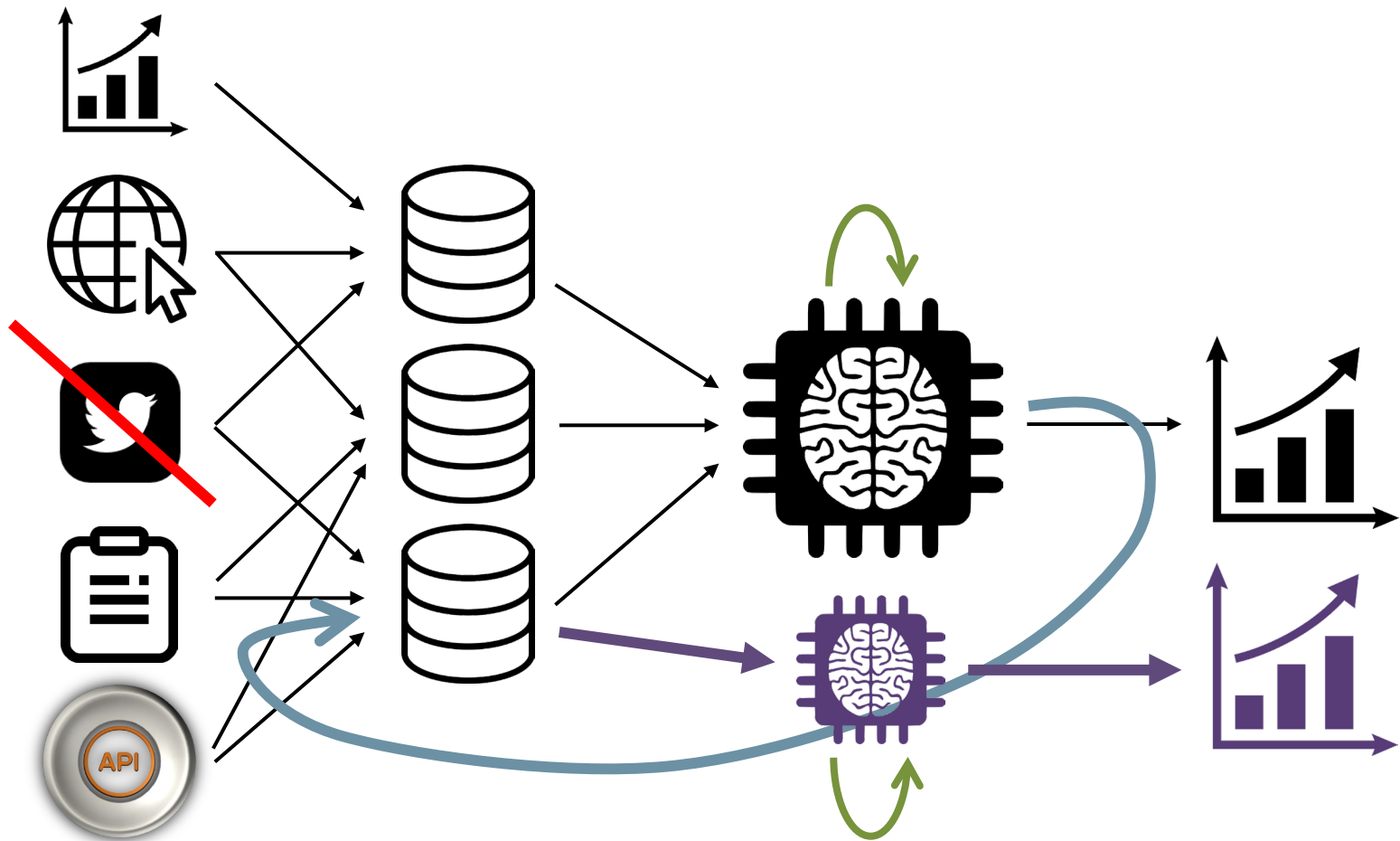# Changing Data Sources
## in the Age of Data Science
## for Official Statistics

Cedric De Boom & Michael Reusens

Statistics Flanders – Belgium

UNECE Machine Learning for Official Statistics

June 6, 2023

# From data source to statistic

STATISTIEK
VLAANDEREN

Vlaamse
overheid

# Benefits of external data sources

| | |
|---|---|
| Broad-spectrum | Covers a wide variety of topics |
| Diversity | A large variety of sources to cover different perspectives |
| Availability | Lots of data is freely and easily accessible |
| Size | Some datasets can be enormous, sometimes even complete |
| Structure | Not only tabular data, but also images, video, text, audio, etc. |
| Timeliness | (Near) up-to-date and real-time information |
| Frequency | Raw data on various, even very fine-grained time scales |
| Granularity | Raw data on various, even fine-grained levels of detail |
| Coverage | Various locations and regions can be filtered and covered |

Vlaamse
overheid

# … but challenges need to be overcome

| | |
|---|---|
| Data quality | Errors, biases, missing values… |
| Data interpretation | Context, meaning, business rules… |
| Data integration | Overcoming diverse structures and formats |
| Selection bias | Ensure representativeness |
| Operationalization bias | Implicit, hidden, and/or production-specific design choices |
| Computational resources | Processing and analyzing large amounts of data |
| Privacy and security | Anonymization, pseudonymization, access management… |
| Data ethics | Data collection and use should adhere to ethical principles |
| Fairness and justness | Neutral, non-discriminatory |
| Cost | Resources, workforce, data purchases… |

STATISTIEK
**VLAANDEREN**

Vlaamse overheid

# Lack of control is an insidious risk



...but    with great amounts of external data comes
          great powerlessness!

          Risk mitigation strategies
          should be front and center
          in your data science agenda and practices!

Vlaamse overheid

# "Data is the new oil" – C.Humby, 2006



Powerful value!

But also:
  Vulnerability!
  Powerlessness!
  Dependency!
  Lack of control!

Vlaamse
overheid

# Types and causes
# of changing data sources

# Overview

Data types and schemas

Sharing and collection technology

Concept drift

Frequency and interruptions

Ownership and discontinuation

Legal properties

Ethics and public perception

STATISTIEK
**VLAANDEREN**

**Confidential** – do not spread

Vlaamse
overheid

# Data types &schemas

= changes in data formats and structure

**Why?**

Accomodate future changes

Technical debt

Improve storage

Increase retrieval effiency

Business rules

…

**Consequences**

Catch errors?

Undetected?

**Mitigation**

Testing, testing, testing!

Data checks &monitoring

Statistical analyses

STATISTIEK
**VLAANDEREN**

Confidential– do not spread

Vlaamse
overheid

# Sharing& collection technology

= storage, cloud, APIs, scraping, external tools, format …

APIs
Endpoint updates
Security patches
Business strategy
Pricing

Vlaamse
overheid

# Concept drift

= data distribution changes between train and test time

**Why?**

Business logic

Variable meaning

Coverage / frequency

Derived data fields!
  i.e. as result of ML model

**Consequences**

Retraining & reevaluation

**Mitigation**

Statistical tests

Monitoring

Vlaamse
overheid

# Frequency& interruptions

= collection or update rate modifications

**Why?**

Deliberate vs random

Technological challenges

Downtime, failures…

**Consequences**

Can lead to concept drift!

**Mitigation**

Statistical tests

Monitoring

# Ownership & discontinuation

= changes in offering or downright shutdown

Consequences

Legal issues, pricing…

Can trigger any other consequence

Mitigation

Redundancy and diversification

Legal contracts / SLAs

# Legalproperties

= legal changes regarding data collection, storage and use

Why?

Privacy laws

Contractual obligations

Mitigation

Redundancy and diversification

Legal contracts / SLAs

Consequences

Renegotiation

Stop the statistical offering

Airtight data management

STATISTIEK
VLAANDEREN

Vlaamse
overheid

Confidential – do not spread

# Ethics & public perception

Why?
Controversial
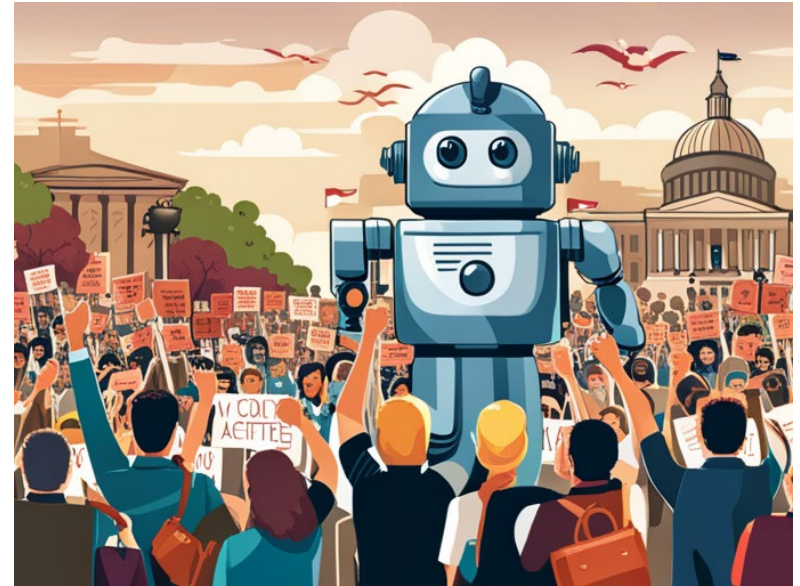Neutrality / bias
Intrusive
Transparency
Accountability
Integrity



Consequences
More scrutiny
More oversight
Patching models
Public trust!
Public policy

STATISTIEK
VLAANDEREN
Confidential – do not spread

Vlaamse
overheid

# In summary

Data types and schemas

Sharing and collection technology

Concept drift

Frequency and interruptions

Ownership and discontinuation

Legal properties

Ethics and public perception

Vlaamse
overheid

# Consequences
## of changing data sources

# Brief overview of consequences

| | |
|---|---|
| Concept Drift | Especially relevant when dealing with long-term trends |
| Model staleness | The model no longer picks up current trends and patterns |
| **Bias and neutrality** | "Garbage in, garbage out" vs neutrality and objectivity |
| Availability | May impact accurate and timely statistics |
| Integration | Beware the domino effect! |
| Extra labor | Take risks into account and allocate resources and time budgets |
| **Breaking changes** | Depending statistics will inevitably change: be transparent! |
| **Quality metrics** | Timeliness, validity, accuracy, completeness, consistency… |

STATISTIEK
VLAANDEREN

Vlaamse
overheid

Confidential – do not spread

# Mitigating
changing data sources

# Mitigation?

Not easy!

Changes are diverse

Consequences are diverse

Required mitigation efforts are time- and resource-consuming

No definite answers...

Highly use-case- and context-dependent

Vlaamse
overheid

# Risk analysis

Identify all potential risks associated with the external data
Use the list in this presentation / paper as a guideline!
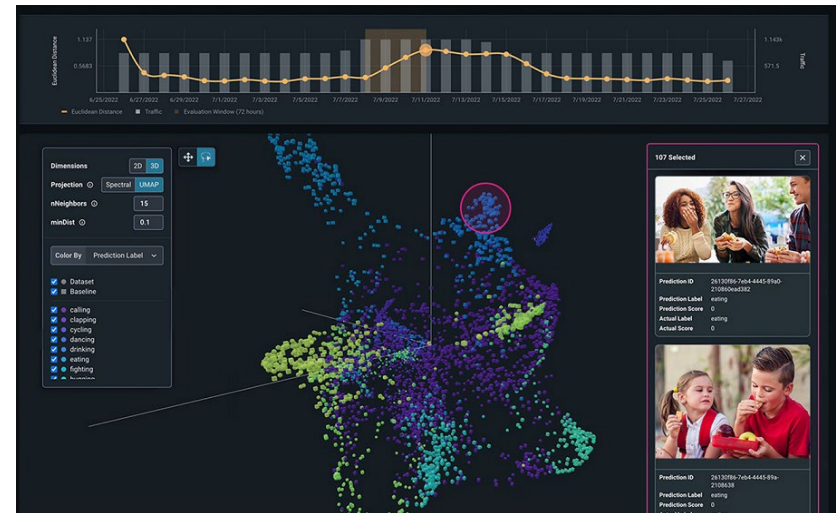


Describe technical
and non-technical aspects

Face the hard truth!

Vlaamse
overheid

# Monitoring

**Monitor everything!**

Record inputs and outputs

Perform statistical tests

Track variables and quantities

**How?**

Reference data sets

Cluster metrics

Visualization and dimensionality reduction

Check predictions against existing domain knowledge

Devise supervised proxy-tasks

# Diversification

Use multiple, redundant data sources if possible

Discrepancies?

Data normalization

Computational overhead

…but is very challenging and not straightforward

# Technicalrobustness

Ensure consistency in the statistical offering

Automated data pipelines

Build resilience against errors, outliers, outages…

Data validation is a part of these pipelines

Thorough unit and integration testing

Failover and deduplication

Security measures

Requires a hefty engineering team,
along with rigorous best practices!

STATISTIEK
**VLAANDEREN**

Vlaamse
overheid

# Legalrobustness

Eliminate unexpected changes, outages and discontinuation

Negotiate tight contracts and SLAs

Specify the legal consequences of non-compliance!

But comes at a significant cost!

# Conclusions

## Risks and consequences

The list is long!

Highly use-case- and context-dependent

This is a story of trade-offs!

But: don't tread lightly on these matters,
especially in the context of official statistics!

## Mitigation strategies

No free lunch

Requires significant resources and a talented workforce

Use this paper and presentation as a guideline / checklist

STATISTIEK
**VLAANDEREN**

Vlaamse
overheid

# Thank you!

**Contact us!**

https://www.vlaanderen.be/
statistiek-vlaanderen-data-science-hub

SCAN ME

Michael Reusens
Data science coordinator
michael.reusens@vlaanderen.be

Cedric De Boom
Senior data scientist
cedric.deboom@vlaanderen.be

STATISTIEK
VLAANDEREN

Vlaamse
overheid