Enabling a Data-informed Public Sector: An Introductory Course to BDTI Essentials

Session 2: Data Cleaning and Transformation

DIGITAL EUROPE PROGRAMME

**Welcome!**
**The webinar will begin shortly.**

**DG CNECT**
Directorate-General for Communications
Networks, Content and Technology

**DG DIGIT**
Directorate-General
for Digital Services

European Commission

# Some housekeeping

- The Chat and Q&A widgets are located to the right of the presenter screen
- Please post session-related questions in Q&A
- Microphones and video are off by default, attendees cannot turn them on
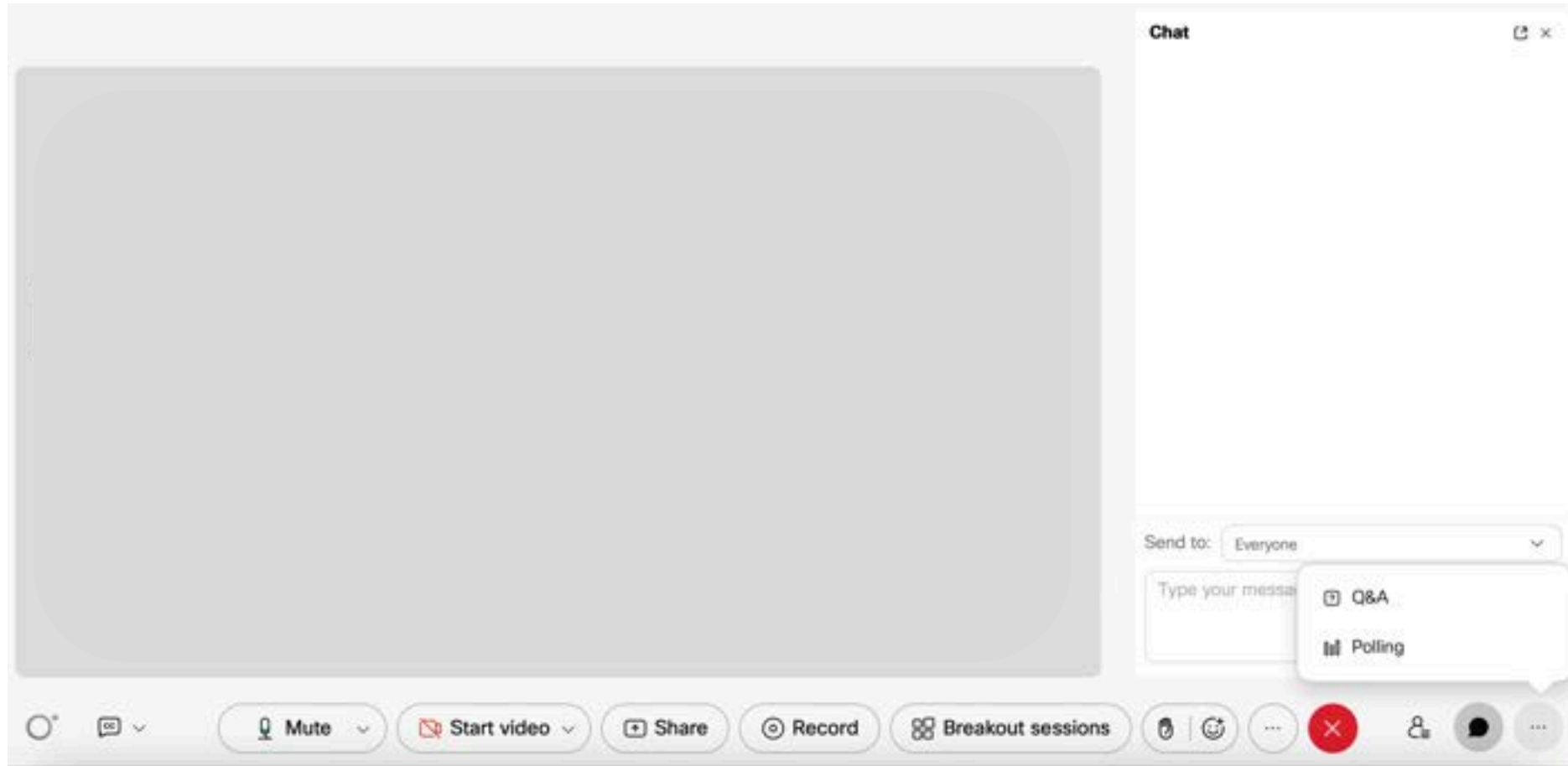- The session will be recorded and available shortly after the end of the session on our website

**Do**
- Participate in the chat during the session
- Ask questions during the Q&A portion at the end
- Give feedback after the session
- Join the dedicated discussion board

**Don't**
- Self-promote
- Disrespect anyone

# Locate Chat and Q&A - Webex

Enabling a Data-informed Public Sector: An Introductory Course to BDTI Essentials

Session 2: Data Cleaning and Transformation

## Agenda

o Welcome [11:00-11:05]
o Session 1 Recap [11.05-11.15]
o Data cleaning and transformation [11:15-12:00]
o Q&A [12:00-12:10]
o Wrap up [12:10 – 12:15]

# Course discussion board

# A bit about you



Geographical Distribution of Webinar Participants



Top Sectors of Participants

# Session 1 –
**BDTI intro,**
**Data Access and Exploration Recap**
**and Exercise review**

# What is the Big Data Test Infrastructure (BDTI) ?

Not **only** for big data, for **public sector in general (open data)**

You have the key ingredients (datasets),
we equip you with the best **open-source tool**
to create amazing recipes for **public good.**

**Six months free of charge service**
for EU public administrations *

**Ready-to-use
data analytics stack** and
support

Cloud platform based on
**open-source** tools

* The cost of the pilot project must fit within the funding boundaries of the BDTI pilot budget

# Open-source tools to support your data journey

**100%**

**Metabase**

**Apache Superset™**

**5. Decision-Making**

**Orchestration**

Apache Airflow

**1. Collection**

**MINIO** — **Data Lake**

**OPENLINK VIRTUOSO UNIVERSAL SERVER**
**mongoDB®** — **Database**

**4. Visualisation**

**Development Environments**

**KNIME**
**jupyter** **H2O.ai**
**R Studio**

**3. Analysis**

**2. Processing**

**elasticsearch**
**Apache Spark™** — **Advanced Processing Engines**

# Who is the Big Data Test Infrastructure (BDTI) for?

**European Public Administrations**
All European Public Administrations at **local, regional and national level** can independently apply for a BDTI pilot project

**Ecosystem with academia and private sector**
Academia, spin-off, startups can apply for pilot projects once there is a **clear collaboration** with a Public Administration which will be the main point of contact for the project (**Master/PhD, GovTech startups)**

**Are you working for a public administration in need of infrastructure for data analytics?**

**https://big-data-test-infrastructure.ec.europa.eu/**

# Use case:
# "Identify Innovative Green Energy Projects from EU Universities"

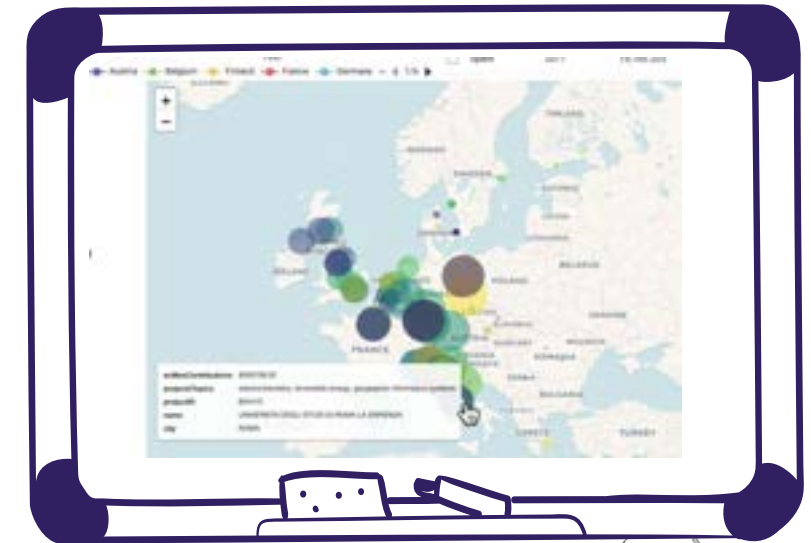Zoi and her team, work in the **Educational Department of a highly polluted EU region.**

Their **purpose** is to identify potential *"green energy partnerships"* between EU universities to address pollution issues.

Zoi reached out the BDTI team to learn how to **connect the dots**, use the free BDTI playground to **create *insights*** and produce a *report*, which can be **easily updated** when new data is available.

They found some interesting **open data:**
- **List of innovation *green energy* projects** across EU countries
- **CO2 emissions** of these countries

# The Data Analysis Process is like cooking a dish

Business challenge/question/problem statement

**Ingredients + tools**

1. **Find** the data you need → *gather the right ingredients, good quality*
2. **Get, clean** and **prepare** your data → *slice and dice*
3. **Analyse** your data → *mix ingredients together and try different combinations*
4. **Present** the results and **create** knowledge → *serve and consume*

**Data→ Information → Knowledge**

# Recap: Zoi Data Journey

**Access the data from multiple sources**

**Clean and transform the data to the correct shape**

**Make sense of data, extract insights and visualise them in a report**

**SESSION 1:**
Data Access and Exploration

**SESSION 2:**
Data Cleaning and Transformation

**SESSION 3:**
Data Blending and Export

**SESSION 4:**
Analytics: aggregation, visualisation, reporting

**SESSION 5:**
Advanced Analytics

**EU Open data portal:**

data.europa.eu

**Formats**:
csv, Excel, JSON

**Tools**:

*Open for Innovation*
**KNIME**

**jupyter**

**R Studio**

**Merge the data and store it in the desired format (end of ETL process)**

**Enhance analytics by gathering data from the internet**

DIGITAL EUROPE PROGRAMME

European Commission

## Learning Resources

In this section you can find some open and online learning resources related to the tools we are going to use in the different sessions

### Jupyter Notebooks

- Phyton Data science Handbook

### R Studio

- R for Data Science

### KNIME

- Extensive Resources for Learning KNIME

- Online self-paced courses

- Explore KNIME Hub for examples

- Documentation to keep your team up to speed with best practices for data science

- KNIME free books

- Ask for help on KNIME Forum

- Videos on KNIMETV

### Statistical learning

- An Introduction to Statistical learning in Python and R

# Exercise Review

# BDTI access refresh

- To have access, public administrations need to *apply on the official website*

- https://big-data-test-infrastructure.ec.europa.eu/apply-bdti_en

- Once you have access follow these steps:

  - Step 1 > Access the portal via the website

  - Step 2 > Navigate to the "Service Catalog" and launch the tools

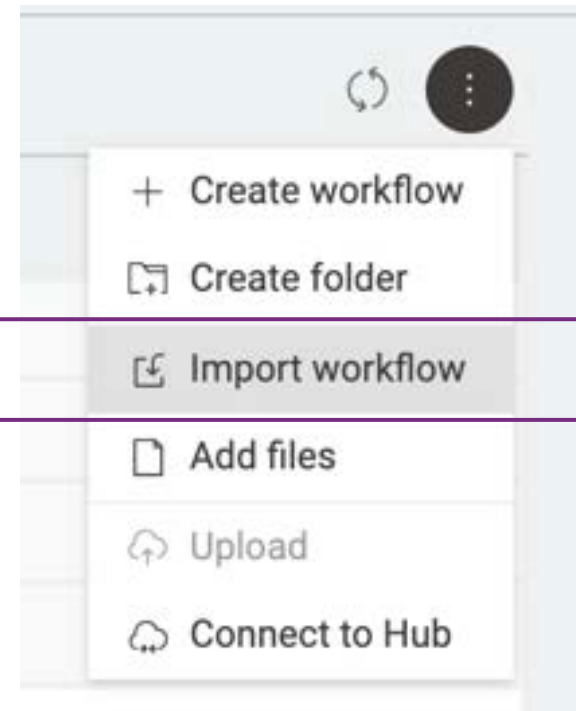  - Step 3 > Go to "My Services" and open the tool you created

# Exercise review

## How to import an existing KNIME workflow?

- Once you are in your working space click on "import workflow"
  - .knwf = 'workflow'
  - .knar = 'archive'

- Remember to reconfigure when needed!

# Exercise review



## Session 1 - Solution

### Data Access and Exploration

Zoi downloads the data from data.europa.eu and from Our World in Data. Your first task is to start building the workflow by accessing the different data sources.

Now that the data access is established, Zoi needs to have a look of the data itself and decide what information is useful to complete her task.

### Access the Data from File(s) and explore the data

1. Read Organization.csv with the **CSV Reader** node
2. Read Project.json with **JSON Reader** node and parse it to a table format with **JSON Path** and **Ungroup** nodes
3. Read The European Science Vocabulary (EuroSciVoc) euroSciVoc.xlsx with the **Excel Reader** node
4. Read the CO2 dataset with the **CSV Reader** node
5. Read the official European Names data with the **Excel Reader** node

**Session 2:**
**Data Cleaning and Transformation**

# Session 2: Data Cleaning and Transformation

| Session | Topic | Duration |
|---------|-------|----------|
| Webinar 1 | Data Access and Exploration | 45 min |
| **Webinar 2** | **Data Cleaning and Transformation** | **75 min** |
| Webinar 3 | Data Blending and Storage | 75 min |
| Webinar 4 | Analytics | 75 min |
| Webinar 5 | Advanced Module: Gathering Data from the Web and Geo Visualisation | 75 min |

- Slides, hands-on exercises, and solutions will be provided on the webinar
- Gitlab repository https://code.europa.eu/bdti/bdti-essentials-course

# Recap of objectives

| Objective | Level |
|---|---|
| **Visualise** total **Energy** projects | EU Countries |
| | Organisations in EU Countries |
| Identify and **Visualise** CO2 emissions | EU Countries |
| **Visualise** energy projects **per year** | EU Countries |
| Explore the Relationship between CO2 emissions and Energy Project Funding | EU Countries |
| **Map** Visualisation of funded projects | EU Universities |

# Outputs according to Zoi goals

**The report/dashboard needs to be informative!**

- Insightful plots

- Tables with a rank of countries according to CO2 emissions

- Bar chart with CO2 emissions

- Line plot with both Horizon 2020 and 2021

- Locate on a map the Universities that got funding for Energy projects from Horizon

# Session 2: Learning outcomes

**At the end of this session, you should be able to:**

- Sign in into the BDTI framework and initialise KNIME.

- Prepare the data for analysis by **cleaning** the dataset.

- Prepare the data for subsequent analysis by **transforming** the dataset to the correct format.

- Address data quality issues.

# Session 2: Data Cleaning and Transformation

- The **core** of the data pipeline

- Why clean and transform?
  - Data are **rarely** clean
  - **Remove** not useful data
  - Remove **repeated** data
  - Handle **missing values**
- **Extract** more information
  - Define rules and mathematical operations
  - Transform at cell or row **level**
- Transform the data to the desired **shape**
  - Rename and resort table columns
  - Split and merge columns
  - Convert data types

**Use Case**

With an overview of the data, Zoi realises that it is not in the best shape to do an analysis, some columns are in the wrong format and she is also missing data.

In this session, she will learn techniques to clean and transform the data.

| Data Access | Data Cleaning and Transformation | Data Blending and Export | Analytics: aggregation, visualisation, reporting | Advanced Analytics |

# Data Cleaning

# Data Cleaning

**Why?**

- "Dirty data"
  - ☐ Unnecessary rows or columns

| Item | Cost | Store Type | City | VAT | Delivery code |
|------|------|------------|------|-----|---------------|
| Table | 50 | Onsite | NY | 14% | XY |
| Couch | 100 | Online | NY | 14% | XX |
| Couch | 100 | Online | NY | 14% | XX |
| Cancelled | ? | ? | ? | ? | ? |
| TV | 70 | Online | NY | 14% | XY |
| Chair | 30 | Onsite | NY | 14% | XY |
| Table | ? | Online | NY | 14% | XX |

# Data Cleaning

**Why?**

- "Dirty data"
  - ☐ Unnecessary rows or columns
  - ☐ Missing or null values

| Item | Cost | Store Type | City | VAT | Delivery code |
|------|------|-----------|------|-----|---------------|
| Table | 50 | Onsite | NY | 14% | XY |
| Couch | 100 | Online | NY | 14% | XX |
| Couch | 100 | Online | NY | 14% | XX |
| Cancelled | ? | ? | ? | ? | ? |
| TV | 70 | Online | NY | 14% | XY |
| Chair | 30 | Onsite | NY | 14% | XY |
| Table | ? | Online | NY | 14% | XX |

# Data Cleaning

**Why?**

- "Dirty data"
  - ☐ Unnecessary rows or columns
  - ☐ Missing or null values
  - ☐ Duplicates

| Item | Cost | Store Type | City | VAT | Delivery code |
|------|------|-----------|------|-----|---------------|
| Table | 50 | Onsite | NY | 14% | XY |
| Couch | 100 | Online | N-Y | 14% | XX |
| Couch | 100 | Online | N-Y | 14% | XX |
| Cancelled | ? | ? | ? | ? | ? |
| TV | 70 | Online | NY | 14% | XY |
| Chair | 30 | Onsite | NY | 14% | XY |
| Table | ? | Online | NY | 14% | XX |

# Data Cleaning

**Why?**

- Keep only what is essential

| Item | Cost | Store Type | City |
|------|------|------------|------|
| Table | 50 | Onsite | NY |
| Couch | 100 | Online | NY |
| Couch | 100 | Online | NY |
| TV | 70 | Online | NY |
| Chair | 30 | Onsite | NY |
| Table | 50 | Online | NY |

# Data Cleaning

**Step 1: Explore data**

Step 2: Identify which columns / rows you need and filter or sort

Step 3: Identify if there are missing values or duplicates you need to deal with

# Step 1 - Explore Data

- Look at raw data (R-studio, Jupyter Notebook, KNIME)
- Review statics view (KNIME)

# Data Cleaning

Step 1: Explore data

**Step 2: Identify which columns / rows you need, and filter or sort**

Step 3: Identify if there are missing values or duplicates you need to deal with

# Step 2 - **Identify** columns

- Horizon Data "Organisations"- *Visualise Energy projects per EU Country and for Organisations per Country*

| RowID | projectID | projectAc... | organisat... | vatNumber | name | shortName | SME | activityTy... | street | postCode | city | country |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Number (inte... | String | Number (inte... | String | String | String | String | String | String | String | String | String |
| Row0 | 879926 | EEN SACHSEN | 999913831 | | INDUSTRIE-U... | IHK C | false | PUB | STRASSE DE... | 09111 | CHEMNITZ | DE |
| Row1 | 879926 | EEN SACHSEN | 994556036 | DE176093942 | ZTS-ZENTRU... | ZTS | true | OTH | INDUSTRIEST... | 01612 | GLAUBITZ | DE |
| Row2 | 879926 | EEN SACHSEN | 999797625 | | INDUSTRIE-U... | IHK DRESDEN | false | PUB | LANGER WEG... | 01239 | DRESDEN | DE |
| Row3 | 879926 | EEN SACHSEN | 999940409 | | | HWK DRESDE... | false | PUB | AM LAGERPL... | | DRESDEN | DE |
| Row4 | 879926 | EEN SACHSEN | 999877844 | DE140857609 | | | false | HES | STRASSE DE... | | Chemnitz | DE |
| Row5 | 879926 | EEN SACHSEN | 999914219 | | INDUSTRIE-U... | IIHK-L | false | PUB | GOERDELERR... | 04109 | LEIPZIG | DE |
| Row6 | 879926 | EEN SACHSEN | 969246408 | | HOCHSCHUL... | Hochschule Z... | false | HES | THEODOR-KO... | 02763 | Zittau | DE |
| Row7 | 879926 | EEN SACHSEN | 999917905 | DE141484057 | AGENTUR FU... | AGIL | false | PRC | LESSINGSTR... | 04109 | LEIPZIG | DE |
| Row8 | 740689 | HEIMDALL | 997710476 | IT015032900... | CENTRO INT... | FONDAZIONE... | false | REC | VIA ARMAND... | 17100 | SAVONA | IT |
| Row9 | 740689 | HEIMDALL | 940057847 | DK37222135 | FREDERIKSB... | FBBR | false | PRC | LOGISMOSE 3 | 3600 | FREDERIKSS... | DK |
| Row... | 740689 | HEIMDALL | 999703438 | ESS0811001G | Departament ... | INT | false | PUB | Carrer Diputa... | 08009 | Barcelona | ES |
| Row... | 740689 | HEIMDALL | 919049102 | IT136697210... | ASSOCIAZIO... | | false | OTH | VIA BERNAR... | 00151 | Roma | IT |
| Row... | 740689 | HEIMDALL | 999778322 | EL094149709 | SPACE HELL... | SPACE HELL... | true | PRC | MESSOGION ... | 153 41 | Aghia Parask... | EL |
| Row... | 740689 | HEIMDALL | 999638739 | ESG62616586 | CENTRE TEC... | CTTC | false | REC | AVINGUDA C... | 08860 | Castelldefels ... | ES |
| Row... | 740689 | HEIMDALL | 996569950 | FR44130005... | UNIVERSITE ... | UNISTRA | false | HES | RUE BLAISE P... | 67081 | Strasbourg | FR |
| Row... | 740689 | HEIMDALL | 935977542 | ESQ0801980D | INSTITUT CA... | ICGC | false | PUB | PARC DE MO... | 08038 | Barcelona | ES |
| Row... | 740689 | HEIMDALL | 997822026 | ESB24352296 | TECNOSYLVA... | | true | PRC | Parque Tecno... | 24009 | Leon | ES |
| Row... | 740689 | HEIMDALL | 999981731 | DE121965658 | DEUTSCHES ... | DLR | false | REC | LINDER HOHE | 51147 | Koln | DE |
| Row... | 740689 | HEIMDALL | 917614763 | GB158535583 | SCOTTISH FI... | SCOTTISH FI... | false | PUB | SCOTTISH FI... | G72 7NA | CAMBUSLANG | UK |
| Row... | 740689 | HEIMDALL | 989367991 | GB830703457 | AVANTI COM... | | true | PRC | COBHAM HO... | EC4V 6EB | London | UK |
| Row... | 740689 | HEIMDALL | 958896217 | ESG55536098 | FUNDACIO D'... | PAU COSTA F... | true | REC | AV. MOSSEN ... | 08552 | Taradell | ES |

- https://data.europa.eu/data/datasets/cordish2020projects?locale=en
- https://data.europa.eu/data/datasets/cordis-eu-research-projects-under-horizon-europe-2021-2027?locale=en

# Step 2 - **Identify** columns

- Horizon Data "Projects" - *Identify Energy **Projects** per Country per **year***

| ids<br>Number (inte... | legalBasi...<br>String | masterCa...<br>String | natures<br>String | objectives<br>String | rcns<br>Number (inte... | startDates<br>String | statuss<br>String | subCalls<br>String | titl...<br>String | topicss<br>String | totalCosts ↓<br>Number (double) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 633053 | H2020-Eurato... | EURATOM-Ad... | | A Roadmap t... | 193159 | 2014-01-01 | SIGNED | EURATOM-Ad... | Imple... | EURATOM | 1,329,689,211.5 |
| 807097 | H2020-EU.3.4. | H2020-IBA-C... | | The challeng... | 218781 | 2014-01-09 | CLOSED | H2020-IBA-C... | Large ... | CS2-GAM... | 247,982.58 |
| 807085 | H2020-EU.3.4. | H2020-IBA-C... | | Engines ITD ... | 216640 | 2014-01-03 | CLOSED | H2020-IBA-C... | Engine... | CS2-GAM... | 228,776,991.75 |
| 945583 | H2020-EU.3.4. | H2020-IBA-C... | | Main objectiv... | 231017 | 2020-01-01 | SIGNED | H2020-IBA-C... | Large ... | CS2-GAM... | 228,776,991.75 |
| 807083 | H2020-EU.3.4. | H2020-IBA-C... | | The Airframe ... | 216645 | 2014-01-02 | CLOSED | H2020-IBA-C... | AIRFR... | CS2-GAM... | 210,183,580.85 |
| 853989 | H2020-EU.3.1. | H2020-JTI-IM... | | The European... | 226660 | 2020-01-01 | SIGNED | H2020-JTI-IM... | EURO... | IMI2-2018. | 207,963,891 |
| 101007873 | H2020-EU.3.1. | H2020-JTI-IM... | | Current anti-t... | 236634 | 2021-06-01 | SIGNED | H2020-JTI-IM... | ACAD... | IMI2-2020. | 185,000,000 |
| 777389 | H2020-EU.3.1. | H2020-JTI-IM... | crisisPrepare... | Paediatric me... | 218657 | 2018-05-01 | SIGNED | H2020-JTI-IM... | conect... | IMI2-2016. | 182,018,216 |
| 737417 | H2020-EU.2.1... | H2020-ECSEL... | | R3-POWERUP... | 210524 | 2017-11-01 | SIGNED | H2020-ECSEL... | 300m... | ECSEL-20... | 180,318,047.5 |
| 681463 | H2020-EU.4.e. | H2020-Adhoc... | | COST is an in... | 198528 | 2015-05-01 | CLOSED | H2020-Adhoc... | COST ... | COST-Net... | 178,169,171 |
| 662338 | H2020-EU.2.1... | ECSEL-2014-2 | | The SeNaTe ... | 198669 | 2015-04-01 | CLOSED | ECSEL-2014-2 | Seven ... | ECSEL-02... | 177,731,624.33 |
| 945542 | H2020-EU.3.4. | H2020-IBA-C... | | The challeng... | 229352 | 2020-01-01 | SIGNED | H2020-IBA-C... | Fast R... | CS2-GAM... | 171,739,896.46 |
| 807081 | H2020-EU.3.4. | H2020-IBA-C... | | The Systems ... | 216639 | 2014-01-01 | CLOSED | H2020-IBA-C... | Syste... | CS2-GAM... | 158,177,923.26 |
| 881603 | H2020-EU.1.2. | H2020-SGA-F... | | This proposal... | 227202 | 2020-04-01 | SIGNED | H2020-SGA-F... | Graph... | SGA-FET-... | 150,000,000 |
| 945539 | | H2020-SGA-F... | | The last of fo... | 229984 | 2020-04-01 | SIGNED | H2020-SGA-F... | Huma... | SGA-FETF... | 150,000,000 |
| 101018100 | H2020-EU.3.1. | H2020-FPA-S... | crisisResponse | Funding of re... | 232261 | 2020-07-24 | SIGNED | H2020-SGA-S... | PROP... | SGA-SC1-... | 150,000,000 |
| 692522 | H2020-EU.2.1... | H2020-ECSEL... | | The TAKE5 pr... | 203403 | 2016-04-01 | CLOSED | H2020-ECSEL... | Techn... | ECSEL-15... | 149,882,181 |
| 101036970 | H2020-EU.3.3. | H2020-LC-GD... | | REFHYNE II w... | 237006 | 2021-10-01 | SIGNED | H2020-LC-GD... | Clean ... | LC-GD-2-2... | 148,956,405 |

- https://data.europa.eu/data/datasets/cordish2020projects?locale=en
- https://data.europa.eu/data/datasets/cordis-eu-research-projects-under-horizon-europe-2021-2027?locale=en

# Step 2 – **Identify** columns

- Horizon Data "EuroSciVoc" - *Identify **Energy** Projects per Country per year*
https://op.europa.eu/en/web/eu-vocabularies/euroscivoc
- Reference vocabulary for the Open Science community and is aligned with Linked Open Data standards.

| RowID | projectID *String* | euroSciVocCode *String* | euroSciVocPath *String* | euroSciVocTitle *String* |
|---|---|---|---|---|
| Row0 | 869855 | /29/101/555/1359 | /social sciences/sociolog... | automation |
| Row1 | 869855 | /25/63/399 | /engineering and technol... | textiles |
| Row2 | 869855 | /25/73/453/58525161 | /engineering and technol... | sensors |
| Row3 | 869855 | /25/73/453/459 | /engineering and technol... | robotics |
| Row4 | 869855 | /29/93/47293307 | /social sciences/psychol... | ergonomics |
| Row5 | 870148 | /25/63/409 | /engineering and technol... | colors |
| Row6 | 883285 | /21/33/121/44109686/5... | /medical and health scien... | ebola |
| Row7 | 883285 | /21/33/137/133/9678651 | /medical and health scien... | pandemics |
| Row8 | 883285 | /21/33/121/44109686/7... | /medical and health scien... | influenza |
| Row9 | 883285 | /29/91/523/1313 | /social sciences/econom... | employment |

- https://data.europa.eu/data/datasets/cordish2020projects?locale=en
- https://data.europa.eu/data/datasets/cordis-eu-research-projects-under-horizon-europe-2021-2027?locale=en

# Step 2 – Identify **columns** and **rows**

- CO2 emissions data: *Identify and Visualise **EU Countries** with the Most and Least **CO2 Emissions***

| | country | year | iso_code | population | gdp | cement_co2 | cement_co2_per_capita | cumulative_cement_co2 | cumulative_co2 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | 1850 | AFG | 3752993.0 | NaN | NaN | NaN | NaN | NaN |
| 1 | Afghanistan | 1851 | AFG | 3767956.0 | NaN | NaN | NaN | NaN | NaN |
| 2 | Afghanistan | 1852 | AFG | 3783940.0 | NaN | NaN | NaN | NaN | NaN |
| 3 | Afghanistan | 1853 | AFG | 3800954.0 | NaN | NaN | NaN | NaN | NaN |
| 4 | Afghanistan | 1854 | AFG | 3818038.0 | NaN | NaN | NaN | NaN | NaN |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 50593 | Zimbabwe | 2017 | ZWE | 14751101.0 | 2.194784e+10 | 0.469 | 0.032 | 18.479 | 753.403 |
| 50594 | Zimbabwe | 2018 | ZWE | 15052191.0 | 2.271535e+10 | 0.558 | 0.037 | 19.037 | 765.199 |
| 50595 | Zimbabwe | 2019 | ZWE | 15354606.0 | NaN | 0.570 | 0.037 | 19.606 | 776.313 |
| 50596 | Zimbabwe | 2020 | ZWE | 15669663.0 | NaN | 0.570 | 0.036 | 20.176 | 786.921 |
| 50597 | Zimbabwe | 2021 | ZWE | 15993525.0 | NaN | 0.570 | 0.036 | 20.745 | 798.217 |

50598 rows × 26 columns

- https://ourworldindata.org/co2-emissions

# Step 2 – Identify **rows**

- EU Country Names data: Filter EU countries (rows)

| # | RowID | AREA String | CODE String | COUNTRY NAME String |
|---|---|---|---|---|
| 29 | Row28 | European Union (EU) | UK | United Kingdom |
| 30 | Row29 | European Free Trade Association (EFTA) | IS | Iceland |
| 31 | Row30 | Eur...tion (EFTA) | LI | Liechtenstein |
| 32 | Row31 | Eur...tion (EFTA) | NO | Norway |
| 33 | Row32 | Eur...tion (EFTA) | CH | Switzerland |
| 34 | Row33 | EU... | TR | Turkey |
| 35 | Row34 | European Neighbourhood Policy (ENP)-South | IL | Israel |
| 36 | Row35 | Other European countries | RU | Russia |
| 37 | Row36 | Non-European countries | AU | Australia |
| 38 | Row37 | Non-European countries | BR | Brazil |

https://ec.europa.eu/eurostat/statistics-explained/images/9/9f/Country_Codes_and_Names.xlsx

# Data Cleaning

Step 1: Explore data

Step 2: Identify which columns/ rows you need, and **filter or sort**

**Step 3: Identify if there are missing values or duplicates you need to deal with**

# Step 3 - Missing Values

- Define how to handle missing values for all columns of a given type
  - Affect all columns that are not explicitly mentioned in the second tab
- Define how to handle missing values for each available column

**Missing Value**



| EEN TOPIC2I... | 919796390 | FR111300216... |
| EEN TOPIC2I... | 952860392 | ⑦ |
| EEN TOPIC2I... | 939950468 | ⑦ |
| EEN TOPIC2I... | 938560167 | ⑦ |
| EEN TOPIC2I... | 999779971 | FR211845001... |
| INNO DK | 902853206 | DK40084711 |
| INNO DK | 901400437 | DK40084606 |
| INNO DK | 902860578 | DK40084746 |
| EEN OUEST H... | 950310359 | FR555136825... |
| EEN OUEST H... | 950310359 | FR555136825... |
| EEN Northern... | 996747848 | NL80144121... |
| EEN-Ukraine | 933428673 | UA36368053 |
| EEN-Ukraine | 938295066 | ⑦ |
| EEN-Ukraine | 927605666 | ⑦ |
| EEN-Ukraine | 899486433 | ⑦ |
| MED4INNO | 915515586 | FR851300224 |

# Step 3 - Duplicates values

- Detect duplicate rows and apply a selected treatment
  - Select columns to check for duplicates
  - Provide options for treating duplicated values

**Duplicate Row Filter**

Flag or Remove Duplicates

Select criteria to keep row

# Data Cleaning: Recap

Step 1: Explore data

Step 2: Filter or sort

Step 3: Missing values and duplicates

| Item | Cost | Store Type | City | VAT | Delivery code |
|------|------|-----------|------|-----|---------------|
| Table | 50 | Onsite | NY | 14% | XY |
| Couch | 100 | Online | N-Y | 14% | XX |
| Couch | 100 | Online | N-Y | 14% | XX |
| Cancelled | ? | ? | ? | ? | ? |
| TV | 70 | Online | NY | 14% | XY |
| Chair | 30 | Onsite | NY | 14% | XY |
| Table | ? | Online | NY | 14% | XX |

| Item | Cost | Store Type | City |
|------|------|-----------|------|
| Table | 50 | Onsite | NY |
| Couch | 100 | Online | NY |
| Couch | 100 | Online | NY |
| TV | 70 | Online | NY |
| Chair | 30 | Onsite | NY |
| Table | 50 | Online | NY |

# Data Transformation

# Data Transformation

- **What** do we transform?
  - Shape of the table (new columns)
  - Content of cells

# Data Transformation

- ## **Why transform data?**

  - Columns in mismatching data type that can lead to problems

  - *What are data types?*
    1. **String**: Text or characters.
    2. **Number**: Numerical values.
    3. **Date & Time**: Points in time or durations.

  - Horizon 2020 column "ID" = number

  - Horizon 2021 column "ID" = string

| RowID | projectID<br>Number (integer) |
|-------|-------------------------------|
| Row0  | 879926 |
| Row1  | 879926 |
| Row2  | 879926 |
| Row3  | 879926 |
| Row4  | 879926 |
| Row5  | 879926 |
| Row6  | 879926 |

| RowID | projectID<br>String |
|-------|---------------------|
| Row0  | 101091483 |
| Row1  | 101091483 |
| Row2  | 101091483 |
| Row3  | 101091483 |
| Row4  | 101091483 |
| Row5  | 101091483 |
| Row6  | 101091483 |

# Data Transformation

- **Why** transform data?
  - Need to extract information that is missing
    - We want the Start Year of the projects
    - "startDate" column, but it comes as a string
    - Need to transform to a date-time format and the extract year!

# Data Transformation

- **Why** transform data?
  - One project can have multiple topic name, so we need to **group** the different topics per project

| RowID | projectID<br>String | euroSciVocTitle<br>String |
|---|---|---|
| Row175 | 852096 | renewable energy |
| Row177 | 852096 | energy conversion |
| Row232 | 951513 | renewable energy |
| Row248 | 759603 | energy conversion |
| Row314 | 742829 | energy and fuels |
| Row536 | 790458 | energy and fuels |
| Row686 | 966685 | solar energy |
| Row800 | 864660 | energy conversion |
| Row811 | 724610 | energy conversion |

| RowID | projectID<br>String | projectsTopics<br>String |
|---|---|---|
| RowID | 852096 | projectsTopics |
| Row2001 | 852096 | renewable energy, energy conversion |

# Data Transformation

- **Why** transform data?
  - One project can have multiple topic name, so we need to **group** the different topics per project

| Project ID | Titles |
|------------|--------|
| 1 | Renewable energy |
| 1 | Energy conversion |
| 2 | solar energy |
| 3 | Solar energy |
| 3 | Energy conversion |
| 4 | Renewable energy |

| Project ID | Concatenate(Titles) |
|------------|---------------------|
| 1 | Renewable energy, Energy conversion |
| 2 | Solar energy |
| 3 | Solar energy, energy conversion |
| 4 | Renewable energy |

Aggregated on Project ID(group) by Concatenation (aggregation method)

# Data Transformation

- **Why** transform data?
  - Some country names are incorrect and can also lead to problems later when merging with all the date
    - Transform string and get them in the correct format

| AREA<br>String | CODE<br>String | COUNTRY NAME<br>String |
|---|---|---|
| European Union (EU) | BE | Belgium |
| European Union (EU) | BG | Bulgaria |
| European Union (EU) | CZ | Czech Republic |
| European Union (EU) | DK | Denmark |
| European Union (EU) | DE | Germany (including former GDR from 1991) |
| European Union (EU) | EE | Estonia |
| European Union (EU) | IE | Ireland |
| European Union (EU) | EL | Greece |
| European Union (EU) | ES | Spain |
| European Union (EU) | FR | France |

| RowID | AREA<br>String | CODE<br>String | COUNTRY NAME<br>String |
|---|---|---|---|
| Row1 | European Union (EU) | BE | Belgium |
| Row2 | European Union (EU) | BG | Bulgaria |
| Row3 | European Union (EU) | CZ | Czech Republic |
| Row4 | European Union (EU) | DK | Denmark |
| Row5 | European Union (EU) | DE | Germany |
| Row6 | European Union (EU) | EE | Estonia |
| Row7 | European Union (EU) | IE | Ireland |
| Row8 | European Union (EU) | EL | Greece |
| Row9 | European Union (EU) | ES | Spain |

# Recap

# Data Transformation: **what** and **why**?

From

raw data

To transformed data in a
format tailored to Zoi's
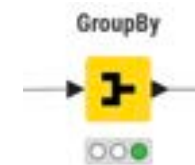needs

# Data Transformation: How?

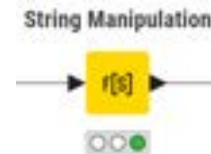- ID Column in Horizon 2020 to the correct data format

- We extracted the Start Year of projects

- Grouped multiple titles into a single row for each project

- Gave Countries the correct wording

# A bit more cleaning

**Sometimes we need to clean again after doing some transformations!**



Clean again!

# Today's Output: Workflow

- All these steps are in one single workflow that we can re-run and do it again!

# Summary

**Now, you should be able to:**

- Sign in into the BDTI framework and initialise KNIME

- Prepare the data for analysis by **cleaning** the dataset

- Prepare the data for subsequent analysis by **transforming** the dataset into the correct format

- Address data quality issues

# Next session

Data Access → Data Cleaning and Transformation → **Data Blending and Export** → Analytics: aggregation, visualisation, reporting → Advanced Analytics

# Practice, practice and practice!

- Slides, hands-on exercises, and solutions will be provided on the webinar Gitlab repository

https://code.europa.eu/bdti/bdti-essentials-course

# Practice, practice and practice!

- Slides, hands-on exercises, and solutions will be provided on the webinar Gitlab repository

https://code.europa.eu/bdti/bdti-essentials-course

## Session 2

### Data Cleaning
With an overview of the data, Zoi asks your help in cleaning the data, since there seems to be duplicates, missing values, and some columns that are not needed.

### Data Transformation
Next, Zoi realises that some columns are in an incorrect data format, so she needs to transform them to be able to use the data in the correct way.

### Clean the data
Clean the **Organisations, Projects and euroSciVoc** data

1. Select columns with the **Column Filter** node.
2. In Organisations 2021, filter the row #45002 using the **Row Filter** node.
3. Use the **Column Resorter** node to set the columns in the order you want.
4. In EuroSciVoc select only energy projects using the **Row Filter** node.

### Transform the data

1. Transform ID column in Horizon2020 to string with the **Number to String** node
2. Transform the startDates and endDates columns format with the **String to Date&Time** node.
3. Extract the Start Year with the **Extract Date&Time Fields** node
4. Rename your new column to "startYear" with the **Column Renamer** node
5. Transform startYear and ID to string with the **Number to String** node
6. Concatenate the different project types in euroSciVoc with the **GroupBy** node
7. Rename this new column with the **Column Renamer** node

**Q&A**

# Your BDTI journey starts here

Congrats! You are on the first step to data-driven innovation.

**BDTI Essentials Course**

**Brainstorm your data project**

**Apply for BDTI Pilot-light process**

# Course discussion board

# Get in touch and follow the BDTI activities

**Are you working for a public administration in need of infrastructure for data analytics?**

**EC-BDTI-PILOTS@ec.europa.eu**

**Visit BDTI's website**

**Subscribe to BDTI's newsletter**

**Subscribe to BDTI's Joinup**

**Register for BDTI Essentials online course**

https://big-data-test-infrastructure.ec.europa.eu/