

# Datos Abiertos en tiempos modernos

**Problemas, soluciones, y algunas ideas.**

@davidgasquez

@davidgasquez



[Create](#)[Home](#)[Competitions](#)[Datasets](#)[Models](#)[Code](#)[Discussions](#)[Learn](#)[More](#)

# Competitions

Grow your data science skills by competing in our exciting competitions. Find help in the [documentation](#) or learn about [Community Competitions](#).

[Host a Competition](#) Search competitions[Filters](#)[All Competitions](#)

Everything, past & present

[Featured](#)

Premier challenges with prizes

[Getting Started](#)

Approachable ML fundamentals

[Research](#)

Scientific and scholarly challenges

[Community](#)

Created by fellow Kagglers

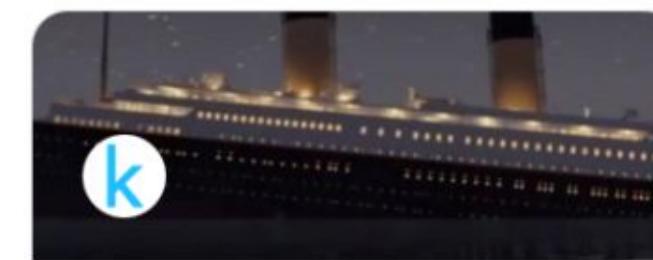
[Playground](#)

Fun practice problems

[Get Started](#)[See all](#)

## New to Kaggle?

These competitions are perfect for newcomers.



### Titanic - Machine Learning from Disaster

Start here! Predict survival on the Titani...

Getting Started

13633 Teams

Knowledge

Ongoing



### House Prices - Advanced Regression Techniques

Predict sales prices and practice feature...

Getting Started

4399 Teams

Knowledge

Ongoing



### Spaceship Titanic

Predict which passengers are transported...

Getting Started

1936 Teams

Knowledge

Ongoing

**train.csv** (61.19 kB)

⬇️ [ ] ➡️

Detail   Compact   Column

11 of 12 columns ▾

**About this file**

contains data

# Survived	# Pclass	▲ Name	▲ Sex	# Age	# SibSp
0	1	891 unique values	male female	65% 35%	 0.42 80
0	3	Braund, Mr. Owen Harris	male	22	1
1	1	Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female	38	1
1	3	Heikkinen, Miss. Laina	female	26	0
1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35	1
0	3	Allen, Mr. William Henry	male	35	0
0	3	Moran, Mr. James	male		0
0	1	McCarthy, Mr. Timothy J	male	54	0
0	3	Palsson, Master. Gosta Leonard	male	2	3
1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27	0
1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14	1
1	3	Sandstrom, Miss. Marguerite Rut	female	4	1
1	1	Bonnell, Miss. Elizabeth	female	58	0
0	3	Saundercock, Mr. William Henry	male	20	0

**test.csv** (28.63 kB)

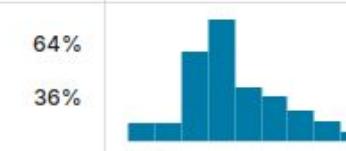
⬇️ [ ] ➡️

Detail   Compact   Column

10 of 11 columns ▾

**About this file**

test data to check the accuracy of the model created

☞ PassengerId	# Pclass	▲ Name	▲ Sex	# Age	# SibSp
892	1309	418 unique values	male female	64% 36%	 0.17 76
892	3	Kelly, Mr. James	male	34.5	0
893	3	Wilkes, Mrs. James (Ellen Needs)	female	47	1
894	2	Myles, Mr. Thomas Francis	male	62	0
895	3	Wirz, Mr. Albert	male	27	0
896	3	Hirvonen, Mrs. Alexander (Helga E Lindqvist)	female	22	1
897	3	Svensson, Mr. Johan Cervin	male	14	0
898	3	Connolly, Miss. Kate	female	30	0
899	2	Caldwell, Mr. Albert Francis	male	26	1
900	3	Abrahim, Mrs. Joseph (Sophie Halaut Easu)	female	18	0
901	3	Davies, Mr. John Samuel	male	21	2
902	3	Ilieff, Mr. Ylio	male		0
903	1	Jones, Mr. Charles Cresson	male	46	0
904	1	Snyder, Mrs. John Pillsbury (Nelle Stevenson)	female	23	1
905	2	Howard, Mr. Benjamin	male	63	1
906	1	Chaffee, Mrs. M. L. Clegg	female	47	1

**train.csv** (61.19 kB)

↓ [ ] >

Detail Compact Column

11 of 12 columns ▾

#### About this file

contains data

#	Survived	#	Pclass	#	Name	#	Sex	#	Age	#	SibSp
0	0	1	1	3	891 unique values		male	65%	0.00 - 0.80 Count: 608		
0	1	3	Braund, Mr. Owen Harris	male	22	1		35%	0.42	80	0
1	1	1	Cumings, Mrs. John Bradley (Florence Thayer)	female	38	1					
1	3	1	Heikkinen, Miss. Laina	female	54	0					
1	1	3	Leino, Mrs. Jacques Heath (Lily May Peel)	female	35	0					
0	0	3	Allen, Mr. William Henry	male	35	0					
0	0	3	Moran, Mr. James	male	54	0					
0	1	1	McCormick, Mr. Timothy J.	male	54	0					
0	3	3	Palsson, Master. Leonid	male	2	3					
1	3	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27	0					
1	2	2	Nasser, Mrs. Nicholas (Adele Embry)	female	14	1					
1	3	3	Sandstrom, Mrs. Marguerite	female	34	0					
1	1	1	Alvarez, Miss. Elizabeth	female	20	0					
0	0	3	Saundercock, Mr. William Henry	male	20	0					

pd.read\_csv()  
model.fit()  
model.predict()

**test.csv** (28.63 kB)

↓ [ ] >

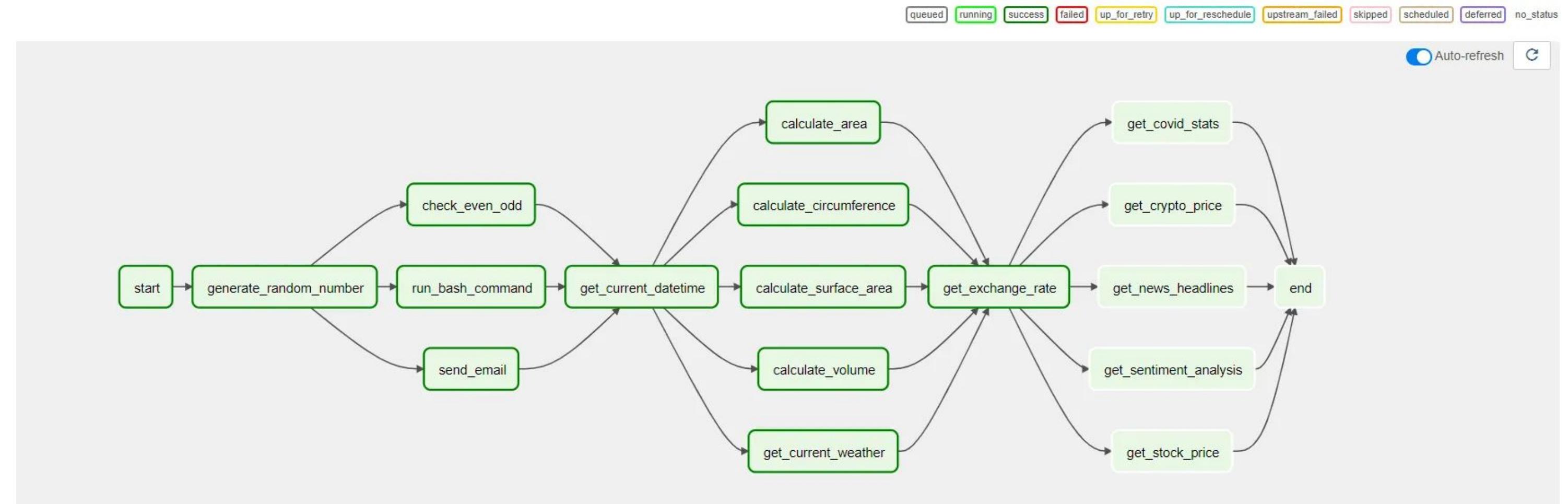
Detail Compact Column

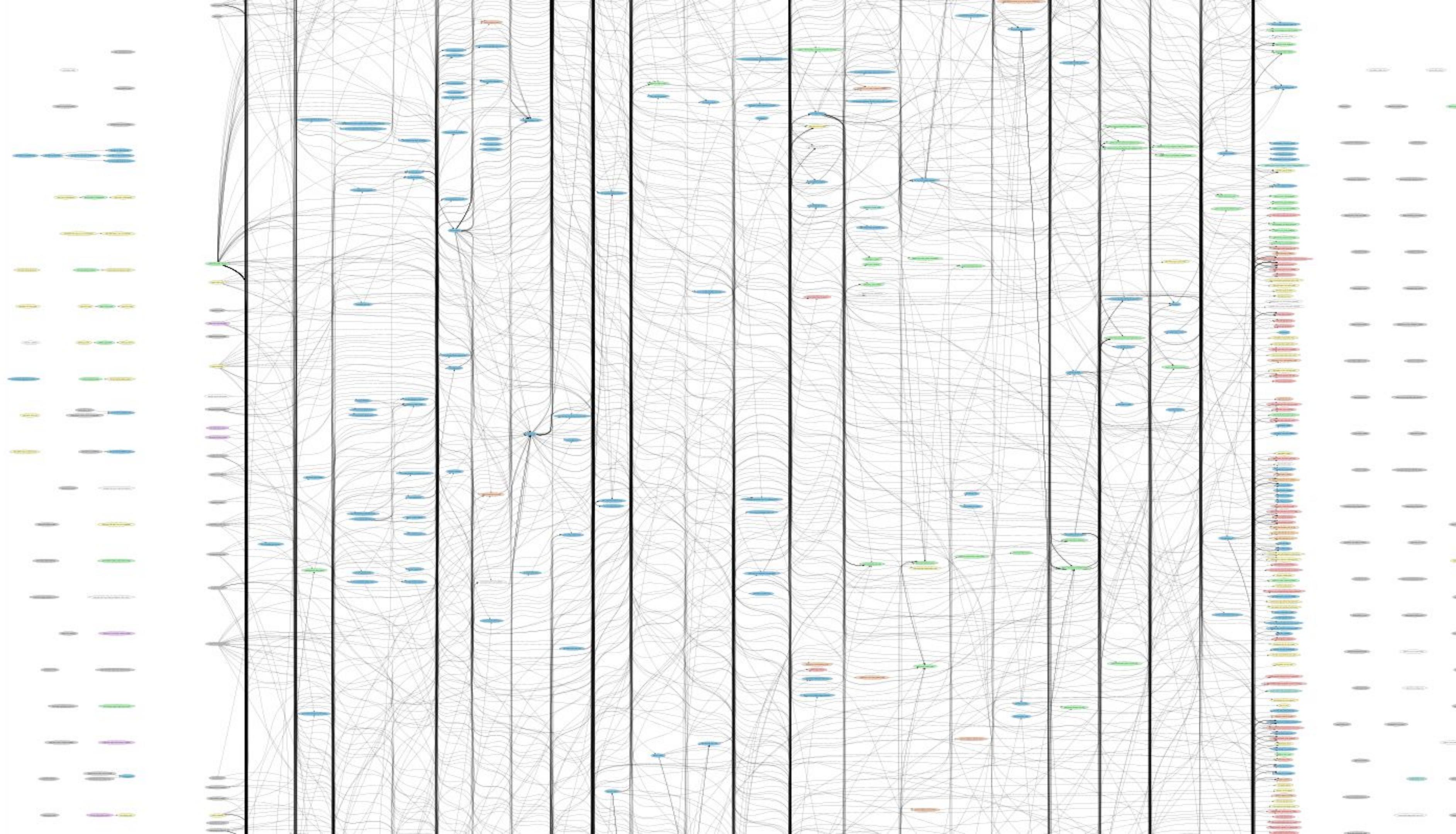
10 of 11 columns ▾

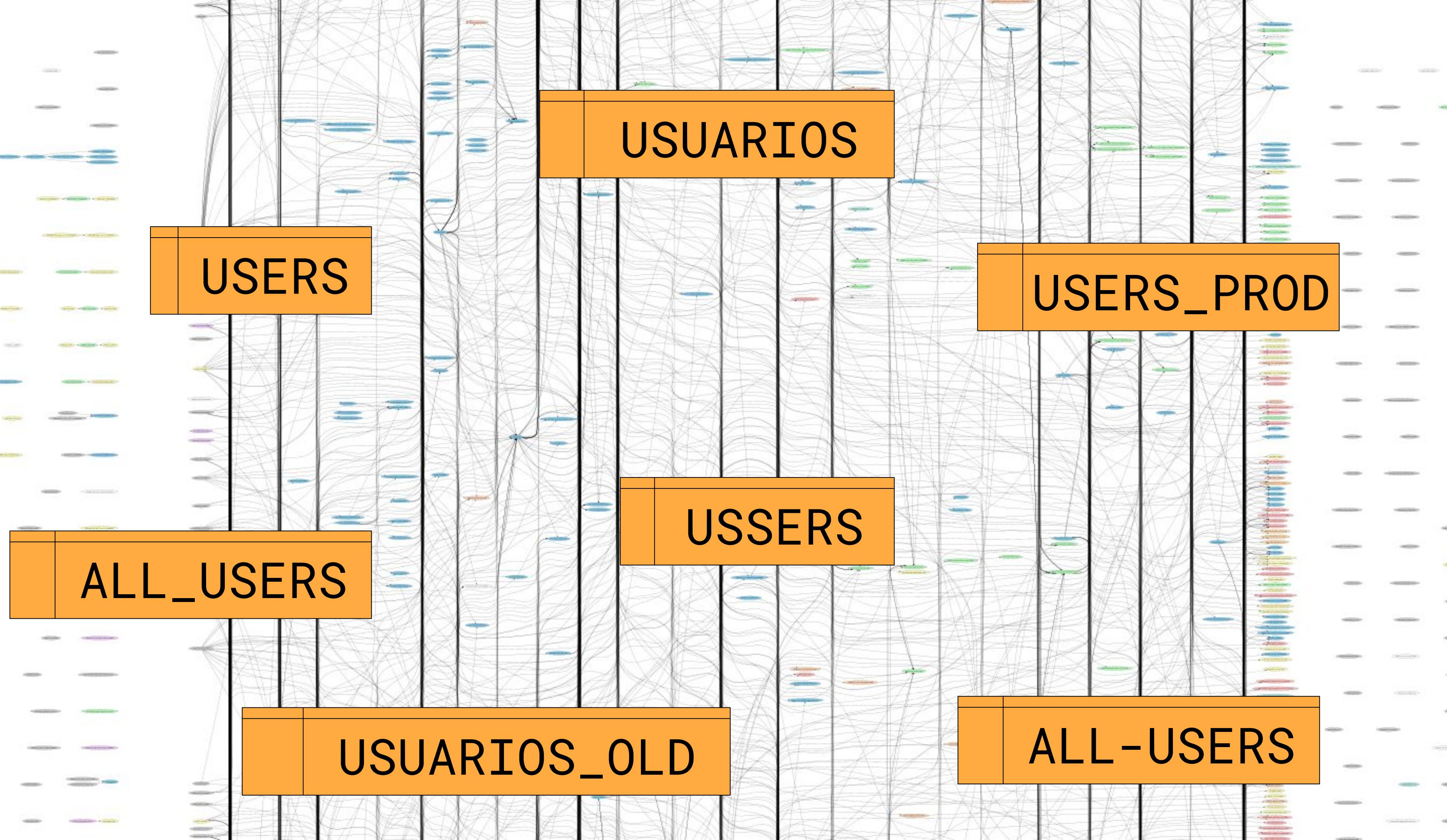
#### About this file

test data to check the accuracy of the model created

#	PassengerId	#	Pclass	#	Name	#	Sex	#	Age	#	SibSp
892	892	1	1	3	418 unique values		male	64%	0.17	76	0
893	893	3			Kelly, Mr. James	male					0
894	894	3			Wilkes, Mrs. James (Ellen Parsons)	female					1
895	895	3			Myles, Mr. Thomas Francis	male					0
896	896	3			Wirz, Mr. Albert	male					0
897	897	3			Hirvonen, Mrs. Alexander (Helga E Lindqvist)	female					1
898	898	3			Svensson, Mr. Johan Cervin	male					0
899	899	3			Connolly, Miss. Kate	female					0
900	900	3			Caldwell, Mr. Albert Francis	male					1
901	901	3			Abrahim, Mrs. Joseph (Sophie Halaut Easu)	female					0
902	902	3			Davies, Mr. John Samuel	male					2
903	903	1			Ilieff, Mr. Ylio	male					0
904	904	1			Jones, Mr. Charles Cresson	male					0
905	905	2			Snyder, Mrs. John Pilkington (Nelle Stevenson)	female					1
906	906	1			Howard, Mr. Benjamin	male					1
					Chaffee, Mrs. Mabel G. Bell	female					1



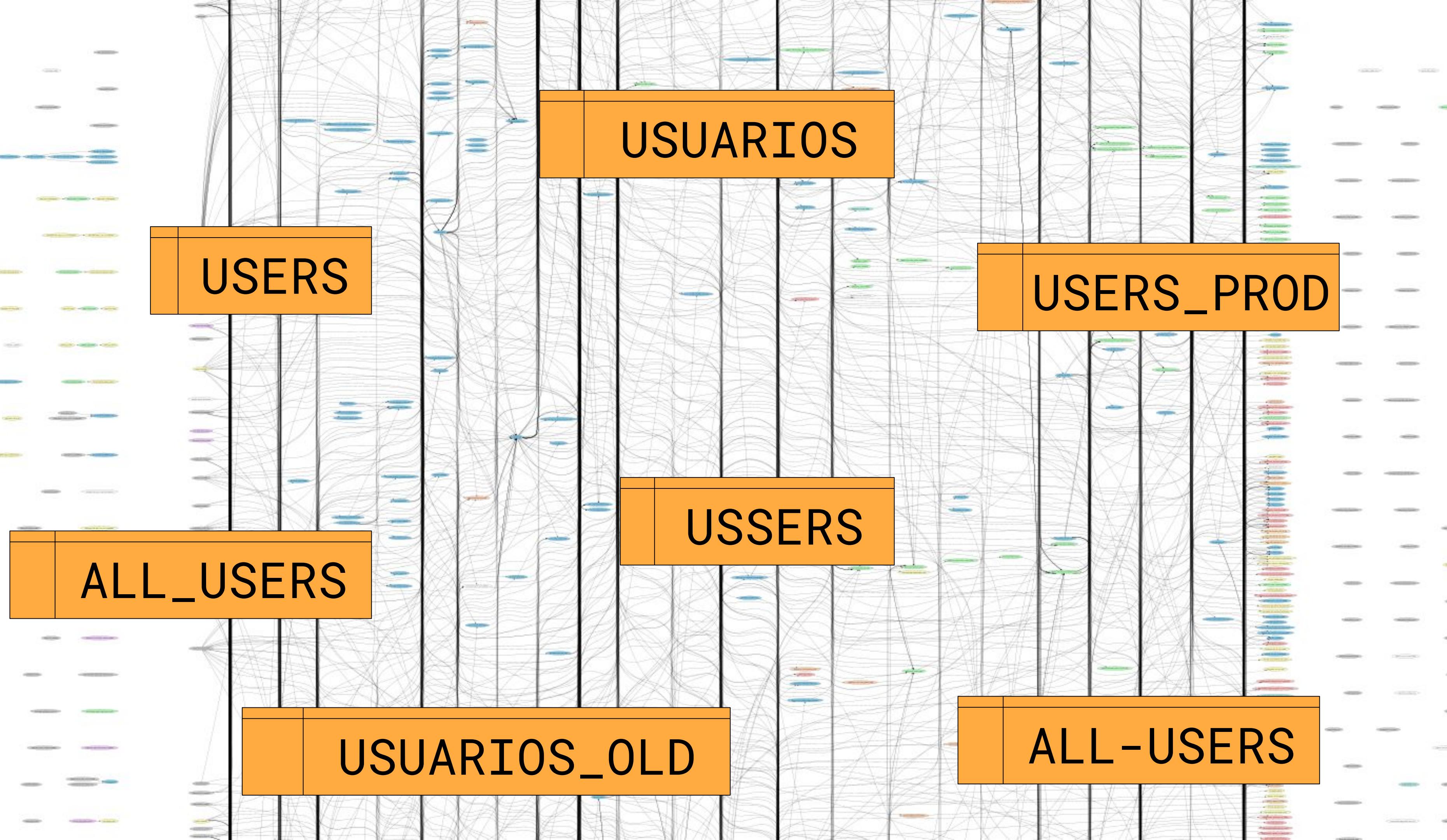






DATA  
SCIENCE

ME



**Inicio**

Definición de URLs

¿Cómo obtener identificadores de objetos utilizando INEbase?

Petición de metadatos

Petición de datos

Glosario Tempus3

Generador de URLs JSON

Generador de gráficos

**API JSON del INE**

El servicio API JSON INE (*Java Script Object Notation*) que se describe en esta sección permite acceder mediante peticiones URL a toda la información disponible en  INEbase, sistema que utiliza el Instituto Nacional de Estadística (INE) para la publicación de la información estadística.

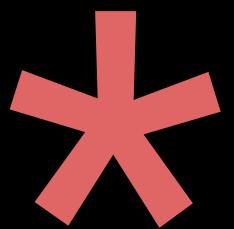
La estructura de las peticiones a través de URL y la simplicidad del formato JSON, hacen que este tipo de aplicaciones sean ampliamente utilizadas para ofrecer datos y metadatos que permiten la explotación automática de la información estadística.

 INEbase es el sistema que utiliza el INE para la publicación de la información estadística. La información que será accesible a través del servicio API JSON del INE que se describe en esta sección, proviene de dos fuentes distintas:

- **Base de datos de difusión (Tempus3).**

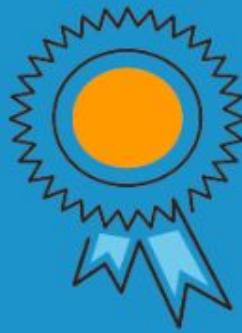


1. Compartir herramientas para trabajar con datos interesantes.
2. Enseñar ejemplos de portales de datos abiertos modernos.



1. Open Data es un ecosistema gigante y diverso.
2. Es más “people problem” que “technical problem”.

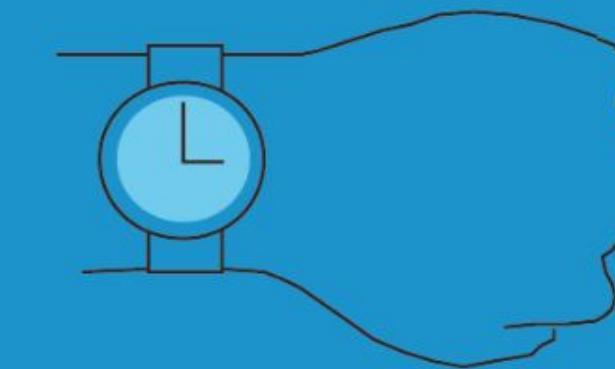
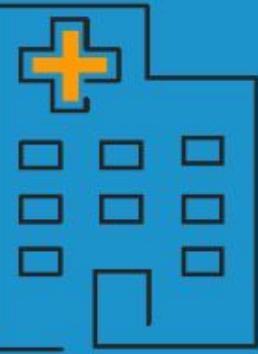
more Open Data can help make  
**better decisions**



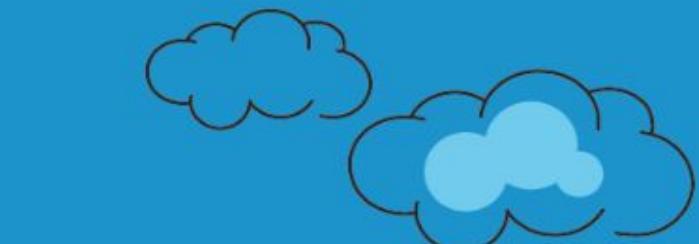
2,549 hours  
**wasted**  
finding parking



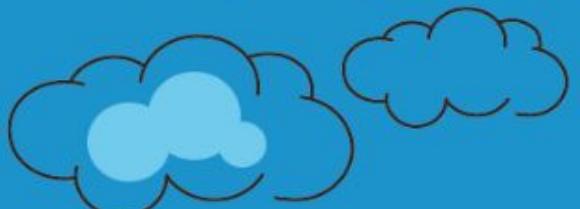
**7,000 lives**  
saved due to  
quicker response



**629 million**  
hours saved is  
equivalent to  
**€ 27.9 bn**



Congestion  
**costs are  
1% of GDP**



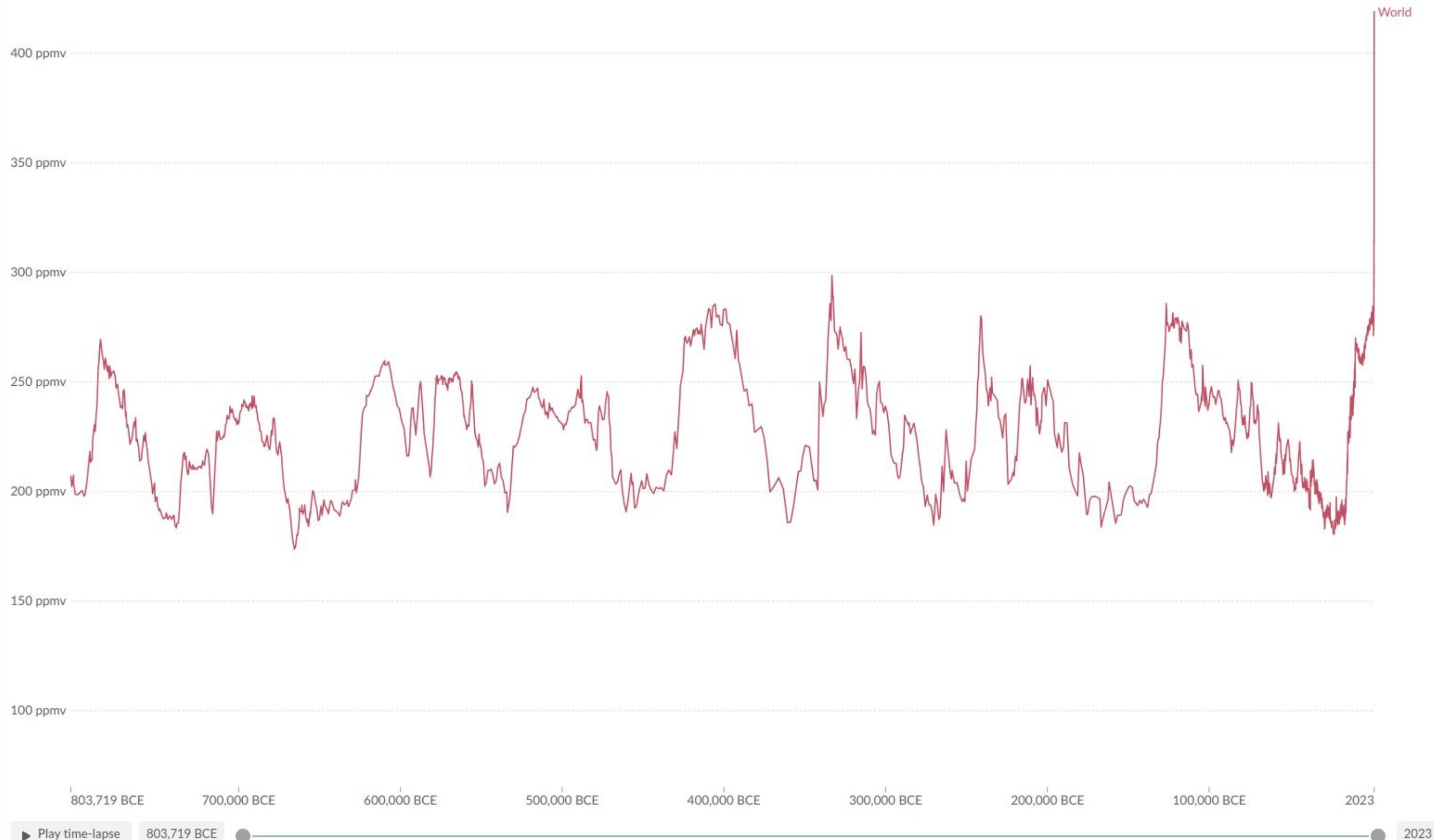
**16% less**  
less energy used



## Carbon dioxide concentrations in the atmosphere

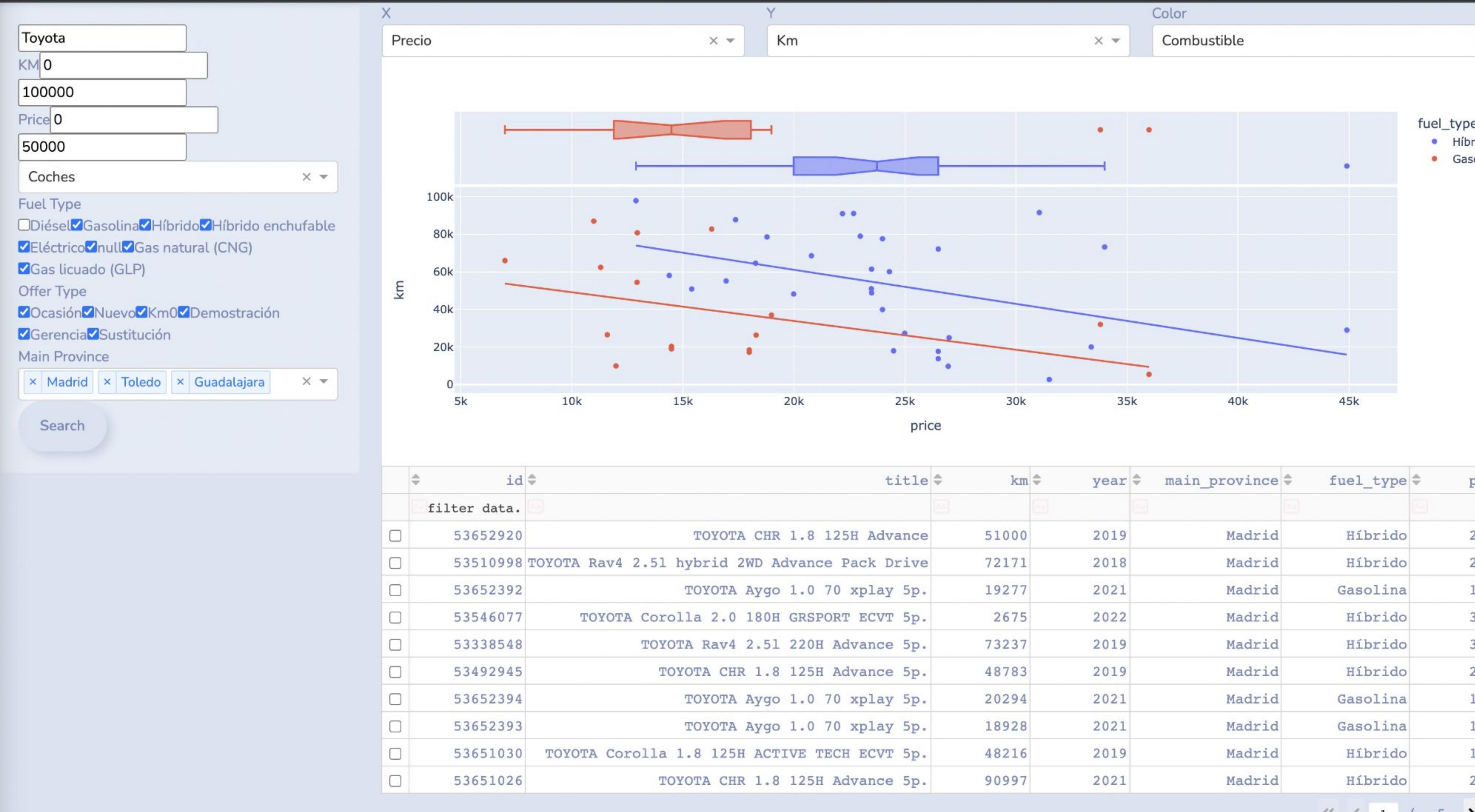
Atmospheric carbon dioxide (CO<sub>2</sub>) concentration is measured in parts per million (ppm). Long-term trends in CO<sub>2</sub> concentrations can be measured at high-resolution using preserved air samples from ice cores.

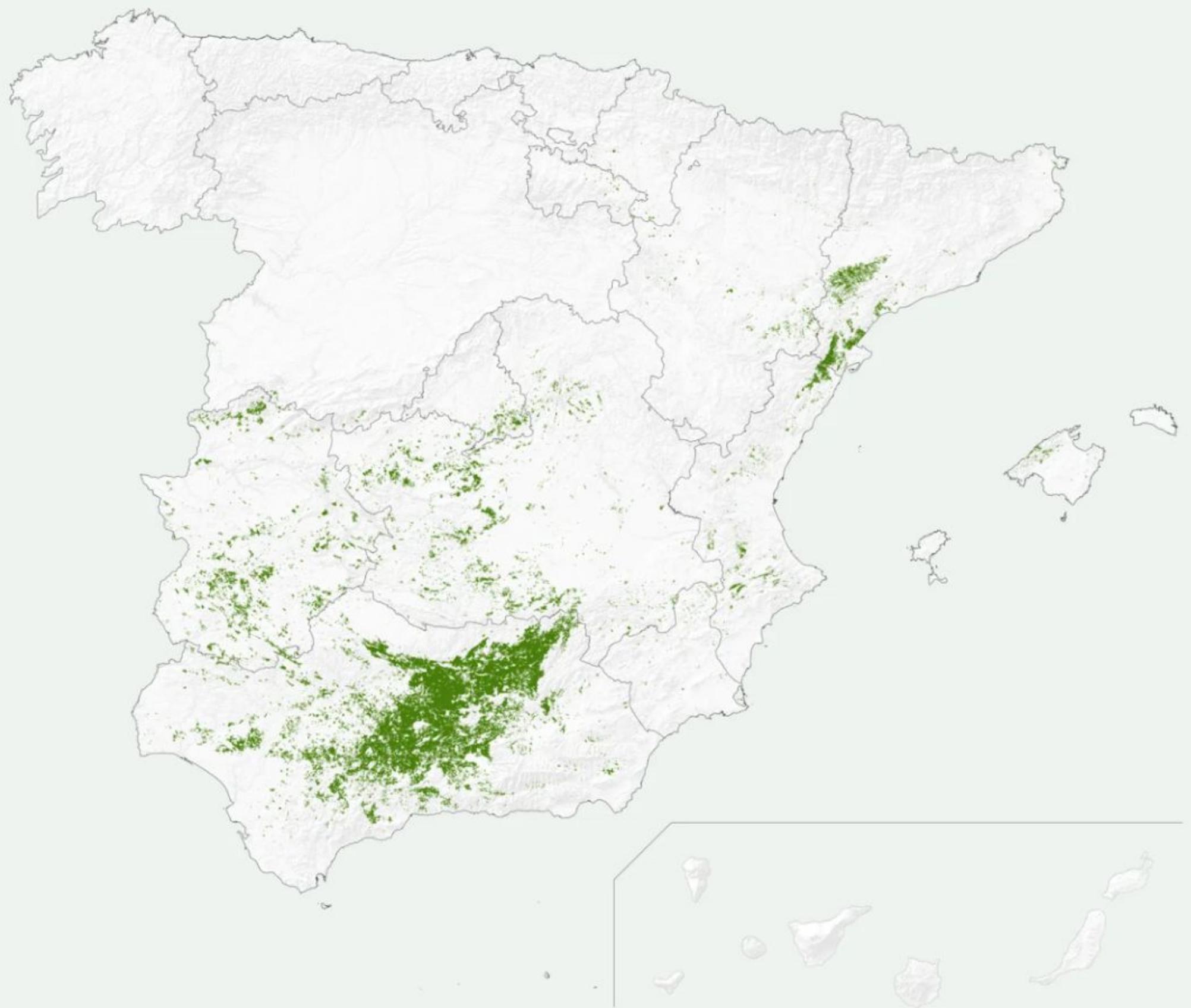
Table Chart



Data source: NOAA Global Monitoring Laboratory - Trends in Atmospheric Carbon Dioxide (2024); EPA based on various sources (2022) – [Learn more about this data](#)  
OurWorldInData.org/climate-change | CC BY

Download  Share  Exit full-screen





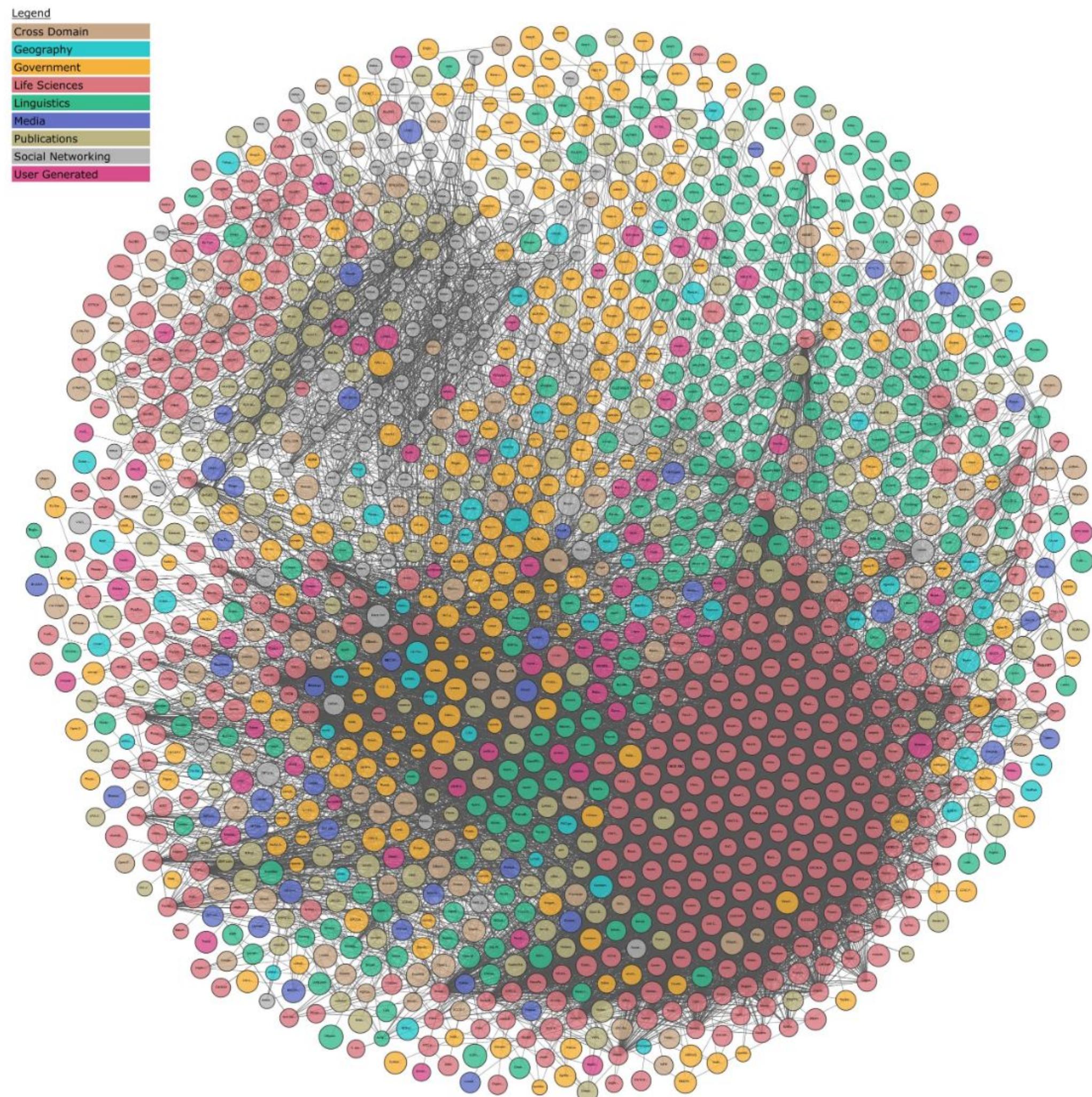
Fuente: [El Orden Mundial](#)



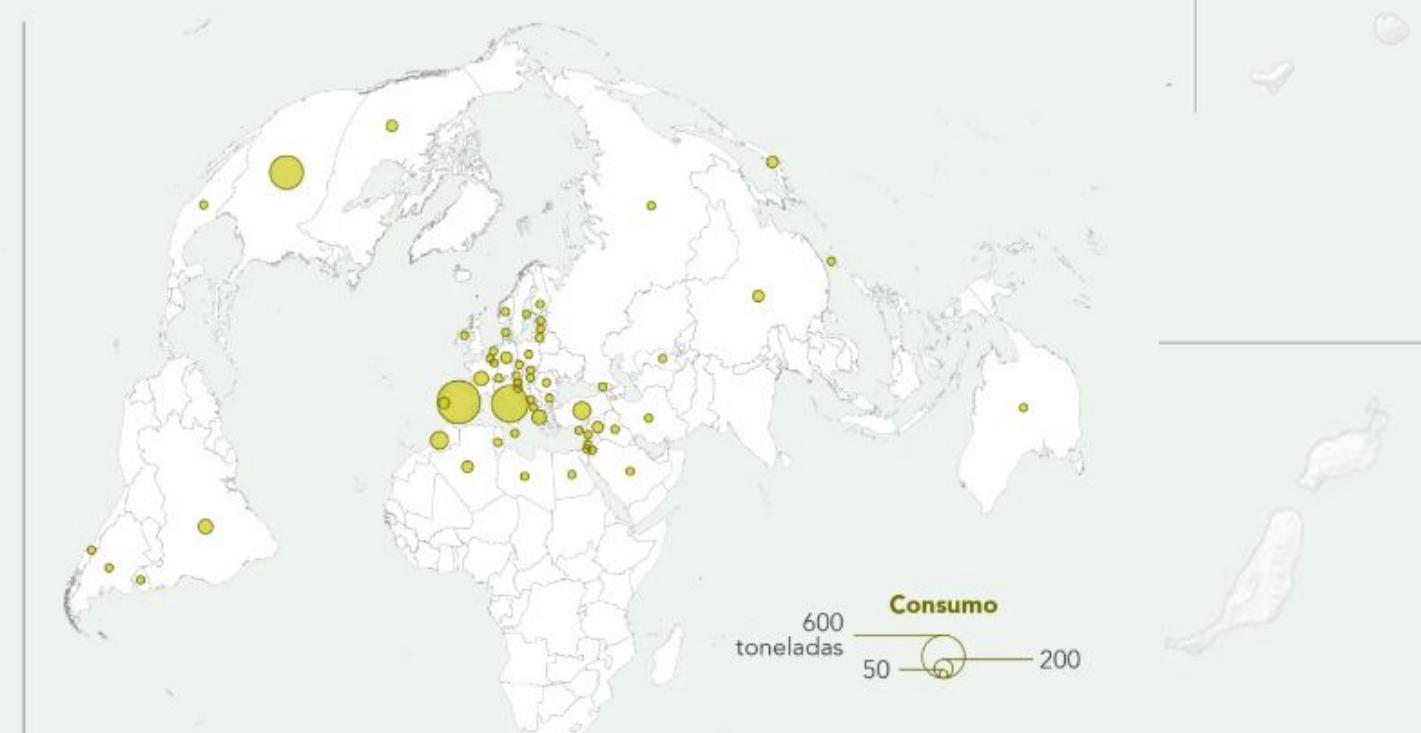
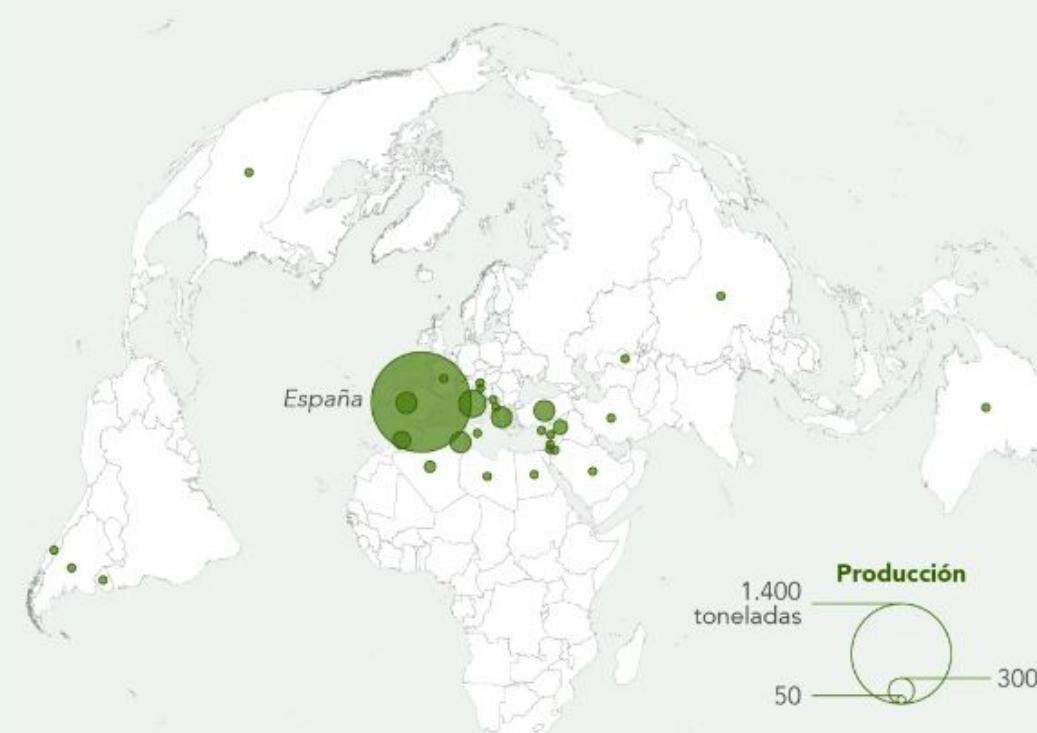
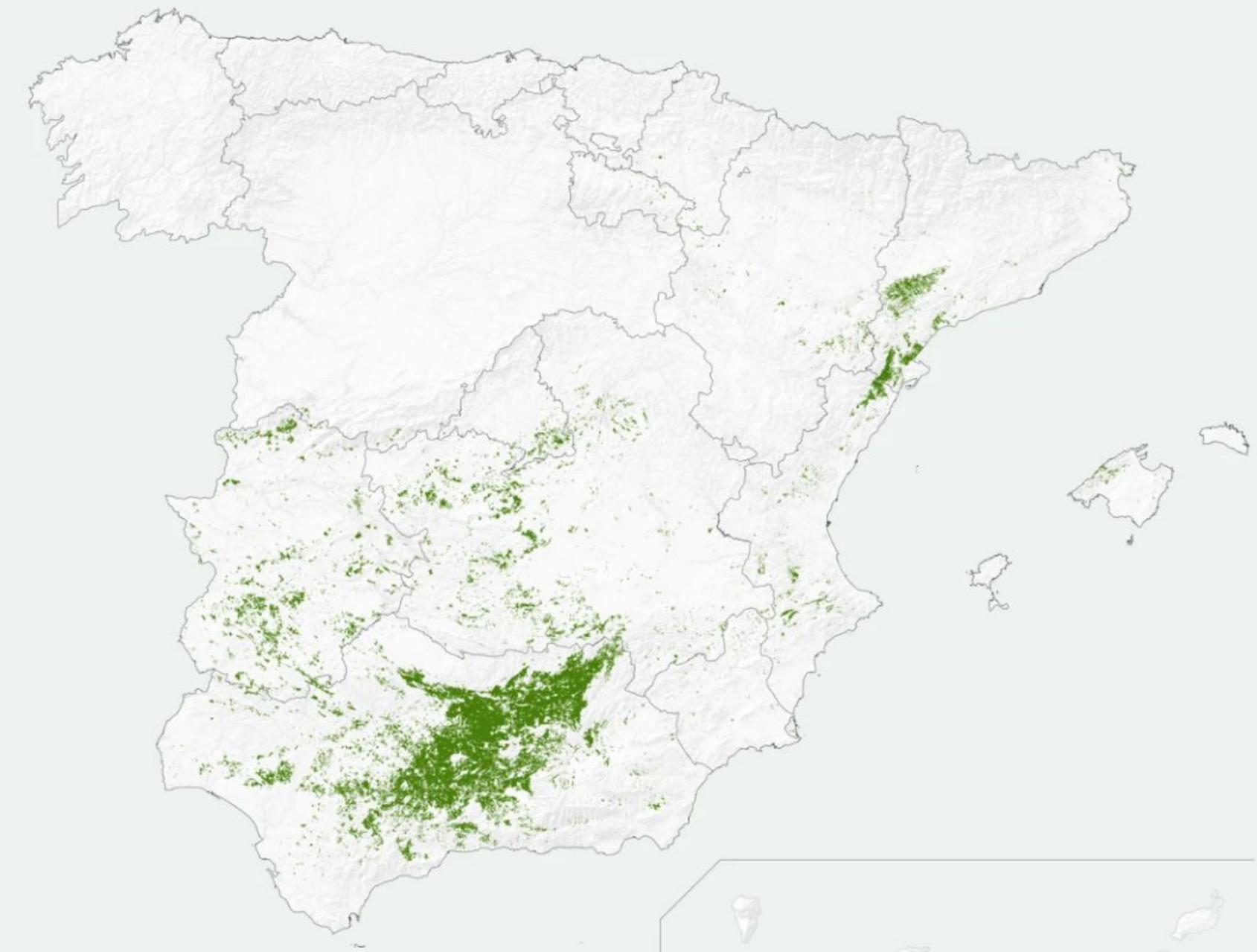
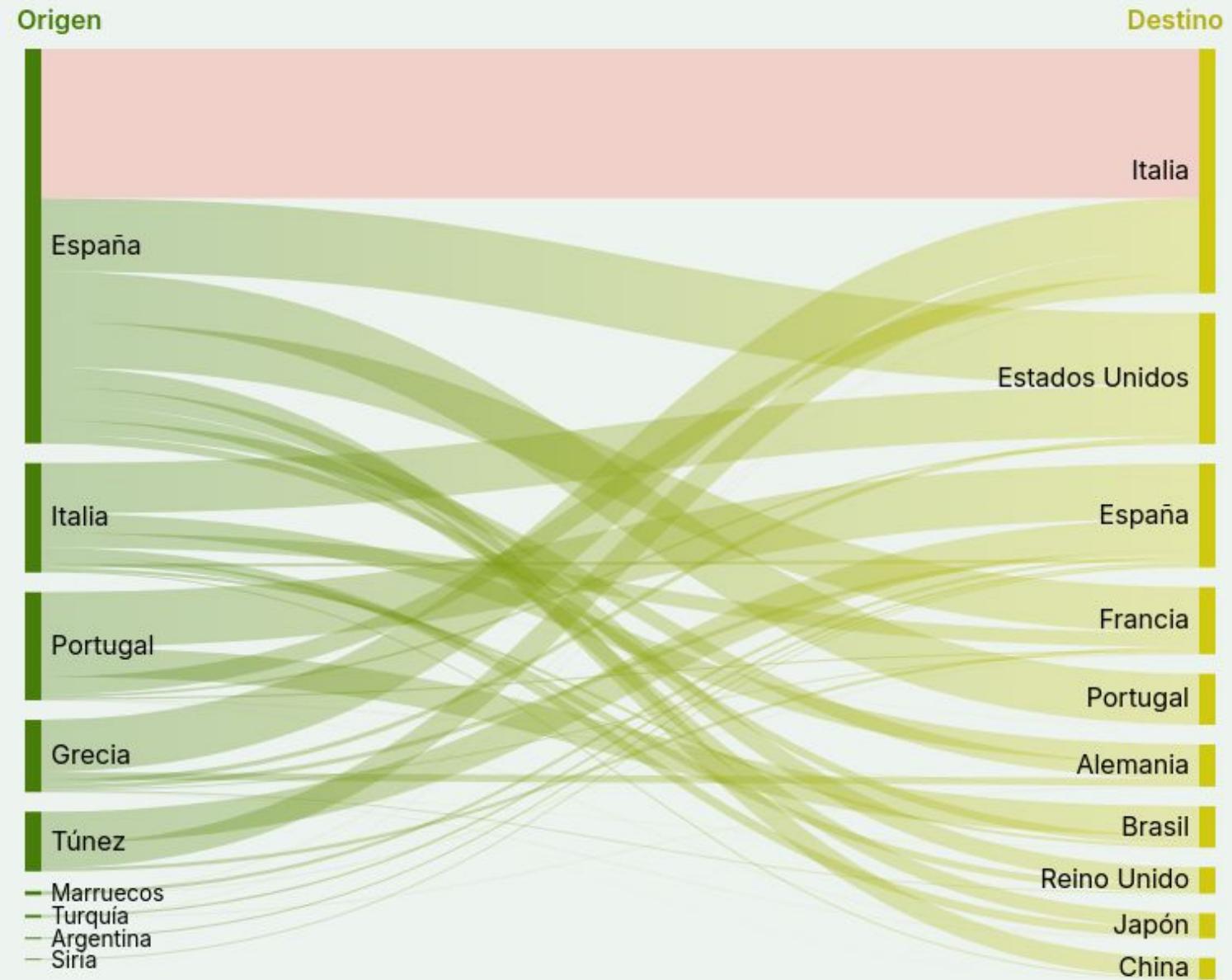
**¡El mundo  
está lleno  
de datos abiertos!**

Oh, oh... Problemas?

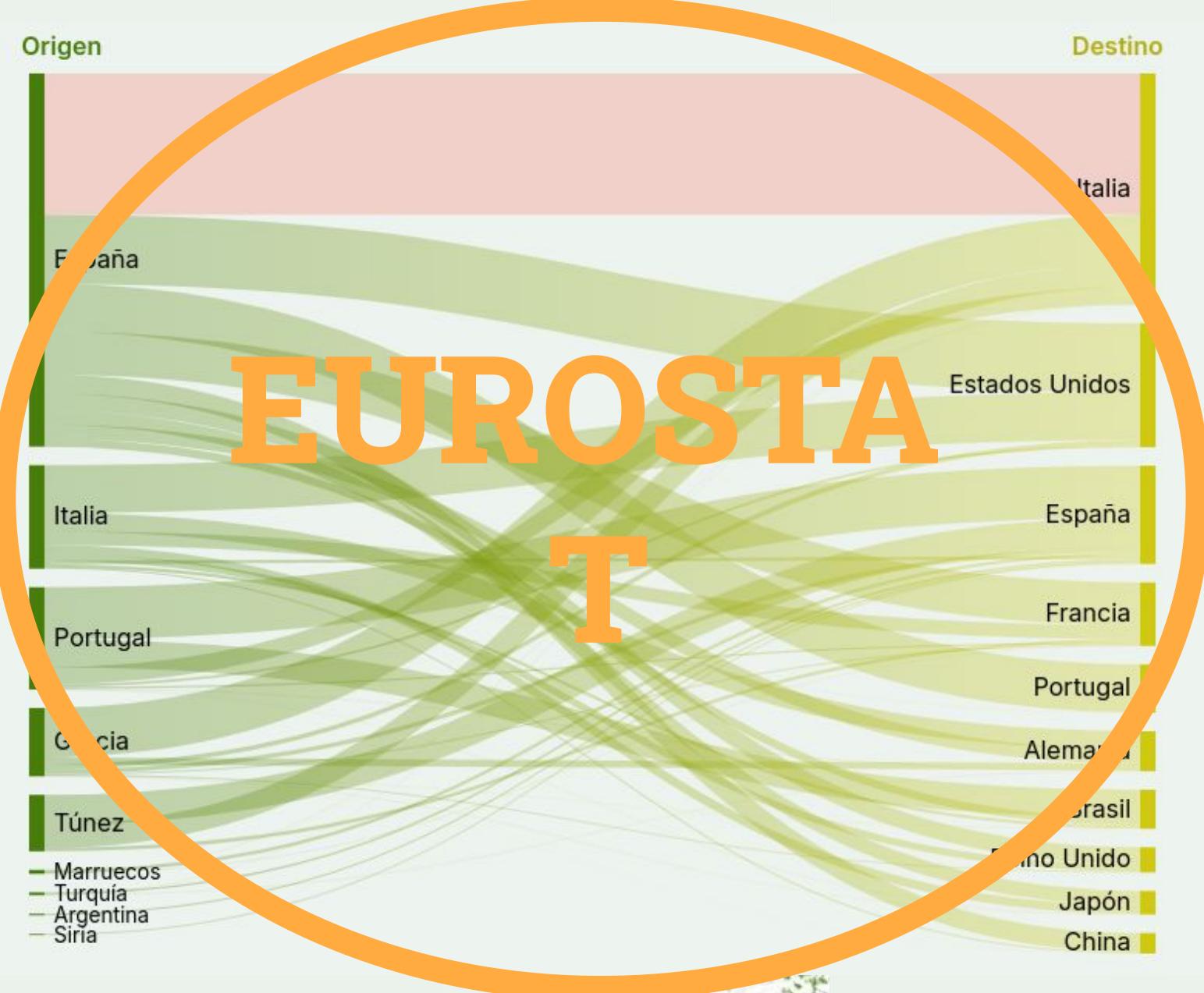
**¡El mundo  
está lleno  
de datos abiertos!**







Fuente: [El Orden Mundial](#)



## Obtener los nombres de los diez organismos que más conjuntos de datos tienen publicados y visualizar el número de éstos

Para realizar esta consulta vamos a tener que agrupar resultados, ordenarlos y limitar el total a 10.

```
select distinct ?label count(?x) as ?num {  
?x a <http://www.w3.org/ns/dcat#Dataset> .  
?x <http://purl.org/dc/terms/publisher> ?publicador.  
?publicador <http://www.w3.org/2004/02/skos/core#prefLabel> ?label.  
}  
group by (?label)  
order by desc(?num)  
limit 10
```

El resultado es este:

label	num
"Gobierno de Aragón"	2659
"Comunidad Autónoma de País Vasco"	2208
"Centro de Investigaciones Sociológicas"	2107
"Ayuntamiento de Málaga"	651
"Ayuntamiento de Gijón"	627
"Xunta de Galicia"	315
"Generalitat Valenciana"	313
"Ayuntamiento de Madrid"	231
"Instituto Nacional de Estadística"	205
"Junta de Castilla y León"	196

## Obtener los nombres de los diez organismos que más conjuntos de datos tienen publicados y visualizar el número de éstos

Para realizar esta consulta vamos a tener que agrupar resultados, ordenarlos y limitar el total a 10.

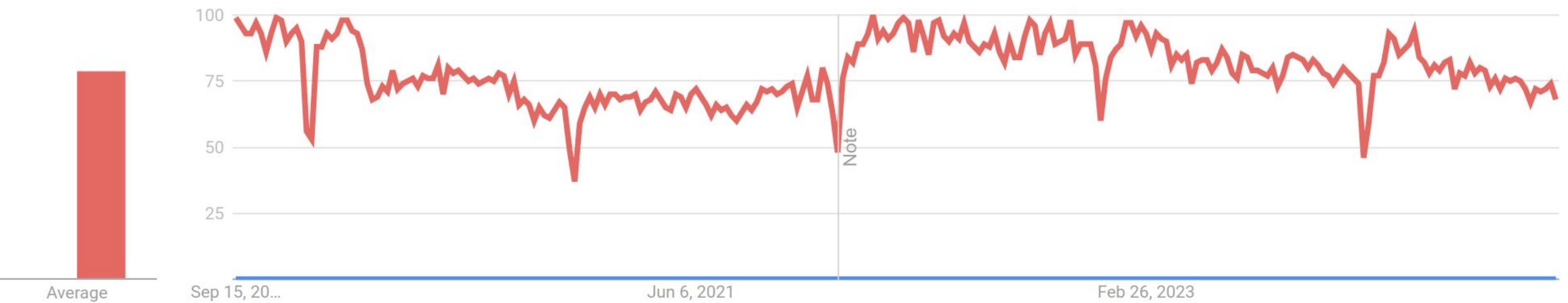
```
select distinct ?label count(?x) as ?num {  
?x a <http://www.w3.org/ns/dcat#Dataset> .  
?x <http://purl.org/dc/terms/publisher> ?publicador.  
?publicador <http://www.w3.org/2004/02/skos/core#prefLabel> ?label.  
}  
group by (?label)  
order by desc(?num)  
limit 10
```

# SPARQL

El resultado es este:

label	num
"Gobierno de Aragón"	2659
"Comunidad Autónoma de País Vasco"	2208
"Centro de Investigaciones Sociológicas"	2107
"Ayuntamiento de Málaga"	651
"Ayuntamiento de Gijón"	627
"Xunta de Galicia"	315
"Generalitat Valenciana"	313
"Ayuntamiento de Madrid"	231
"Instituto Nacional de Estadística"	205
"Junta de Castilla y León"	196

## Interest over time ?



Wikidata Query Service Examples Help More tools Query Builder 文 English

```
1 SELECT
2 ?asteroidLabel
3 ?discovered
4 ?discovererLabel
5 WHERE {
6   ?asteroid wdt:P31 wd:Q3863; # Retrieve instances of "asteroid"
7     wdt:P61 ?discoverer; # Retrieve discoverer of the asteroid
8     wdt:P575 ?discovered; # Retrieve discovered date of the asteroid
9     SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
10
11 BIND (?asteroidLabel AS ?asteroid_label)
12 BIND (?discovererLabel AS ?discoverer_label)
13 }
14 ORDER BY DESC(?discovered)
```

5:13:46 PM GMT+2, Sep 24, 2024  
Data updated a few seconds ago

Table ? 97876 results in 28536 ms </> Code Download Link

Search

asteroidLabel	discovered	discovererLabel
2023 HT4	24 April 2023	Mount Lemmon Survey
2023 FY3	25 March 2023	Catalina Sky Survey
2023 FY3	25 March 2023	Kacper Wierzchoś
2023 DZ2	27 February 2023	<a href="http://www.wikidata.org/.well-known/genid/1fc11aca1fac426d0763053510afd316">http://www.wikidata.org/.well-known/genid/1fc11aca1fac426d0763053510afd316</a>
2023 BU	21 January 2023	Gennady Borisov
2022 YG	15 December 2022	Gennady Borisov
2022 DX	22 February 2022	Erwin Schwab

```

1 SELECT
2 ?asteroidLabel
3 ?discovered
4 ?discovererLabel
5 WHERE {
6   ?asteroid wdt:P31 wd:Q3863; # Retrieve instances of "asteroid"
7     wdt:P61 ?discoverer; # Retrieve discoverer of the asteroid
8     wdt:P575 ?discovered; # Retrieve discovered date of the asteroid
9     SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
10
11 BIND (?asteroidLabel AS ?asteroid_label)
12 BIND (?discovererLabel AS ?discoverer_label)
13 }
14 ORDER BY DESC(?discovered)

```


Table

asteroidLabel	discovered	discovererLabel
2023 HT4	24 April 2023	Mount Lemmon Survey
2023 FY3	25 March 2023	Catalina Sky Survey
2023 FY3	25 March 2023	Kacper Wierzchoś
2023 DZ2	27 February 2023	<a href="http://www.wikidata.org/.well-known/">http://www.wikidata.org/.well-known/</a>
2023 BU	21 January 2023	Gennady Borisov
2022 YG	15 December 2022	Gennady Borisov
2022 DX	22 February 2022	Erwin Schwab

asteroid name or id, discovered date, discoverer

change the column names accordingly to fit WIKIDATAs ones



Sure! Here's a WIKIDATA SPARQL query that retrieves asteroid data including the asteroid name or ID, discovered date, and discoverer:

```

ruby
Copy code

SELECT ?asteroid ?discovered ?discovererLabel
WHERE {
  ?asteroid wdt:P31 wd:Q3863; # Retrieve instances of "asteroid"
    wdt:P61 ?discoverer; # Retrieve discoverer of the asteroid
    wdt:P575 ?discovered. # Retrieve discovered date of the asteroid
  SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
}

```

In the above query, we use the `SELECT` clause to specify the variables we want to retrieve, which are `?asteroid`, `?discovered`, and `?discovererLabel`. The `?asteroid` variable represents the ID or name of the asteroid, the `?discovered` variable represents the discovered date of the asteroid, and the `?discovererLabel` variable represents the discoverer of the asteroid.

In the `WHERE` clause, we use the `wdt:P31` property to filter for instances of "asteroid" (Q3863) in WIKIDATA. We then use the `wdt:P61` property to retrieve the discoverer of the asteroid, and the `wdt:P575` property to retrieve the discovered date of the asteroid. Finally, we use the `SERVICE` clause with `wikibase:label` to retrieve the English label for the discoverer using the `?discovererLabel` variable.

Note: Please note that the results may be limited to the available data in WIKIDATA, and some asteroids may not have complete information.

# Datasets: ⚡ datonic/wikidata\_asteroids

like 0

Modalities: Text

Formats: parquet

Size: 10K - 100K

Libraries: Datasets

pandas

Croissant

+1 License: mit

Dataset card

Viewer

Files and versions

Community 11

Settings

## Dataset Viewer

Auto-converted to Parquet

API

Embed

Full Screen Viewer

Split (1)

train · 97.4k rows

Search this dataset

SQL Console

asteroidLabel  
large\_string · lengths



2023 HT4

discovered  
large\_string · lengths



2023-04-24T00:00:00Z

discovererLabel  
large\_string · lengths



Mount Lemmon Survey

2023 FY3

2023-03-25T00:00:00Z

Catalina Sky Survey

2023 FY3

2023-03-25T00:00:00Z

Kacper Wierzchoś

2023 DZ2

2023-02-27T00:00:00Z

<http://www.wikidata.org/.well-known/genid/1fc11aca1fac426d0763053510afd316>

2023 BU

2023-01-21T00:00:00Z

Gennady Borisov

2022 YG

2022-12-15T00:00:00Z

Gennady Borisov

2022 DX

2022-02-22T00:00:00Z

Erwin Schwab

2022 AE1

2022-01-06T00:00:00Z

Mount Lemmon Survey

< Previous

1

2

3

...

975

Next >

Downloads last month

2

Use this dataset

Edit dataset card

⋮

Size of downloaded dataset files:  
848 kB

Size of the auto-converted Parquet files:  
848 kB

Number of rows:  
97,422



## Índice de Precios de Consumo. Base 2021. Medias anuales

Resultados nacionales

### Índices nacionales de clases

Unidades: Índice, Tasas



#### ▶ Seleccione valores a consultar

Clases	Tipo de dato	Periodo
<input type="text"/> <a href="#">Índice general</a> <a href="#">0111 Pan y cereales</a> <a href="#">0112 Carne</a> <a href="#">0113 Pescado y marisco</a> <a href="#">0114 Leche, queso y huevos</a> <a href="#">0115 Aceites y grasas</a>	<input type="text"/> <a href="#">Media anual</a> <a href="#">Variación de las medias anuales</a>	<input type="text"/> <a href="#">2023</a> <a href="#">2022</a> <a href="#">2021</a> <a href="#">2020</a> <a href="#">2019</a> <a href="#">2018</a>
Seleccionados: 93	Seleccionados: 1	Seleccionados: 1
Total: 93	Total: 2	Total: 22

#### ▶ Elija forma de presentación de la tabla

<b>Clases</b>	

Decimales a mostrar:

#### ▶ Notas

Total: 93 series y 93 datos

Consultar selección

Consultar todo





## Índice de Precios de Consumo. Base 2021. Medias anuales

Resultados nacionales

### Índices nacionales de clases

Unidades: Índice, Tasas



#### ▶ Seleccione valores a consultar

Clases	Tipo de dato	Periodo
<input type="text"/> <a href="#">Índice general</a> <a href="#">0111 Pan y cereales</a> <a href="#">0112 Carne</a> <a href="#">0113 Pescado y marisco</a> <a href="#">0114 Leche, queso y huevos</a> <a href="#">0115 Aceites y grasas</a>	<input type="text"/> <a href="#">Media anual</a> <a href="#">Variación de las medias anuales</a>	<input type="text"/> <a href="#">2023</a> <a href="#">2022</a> <a href="#">2021</a> <a href="#">2020</a> <a href="#">2019</a> <a href="#">2018</a>
Seleccionados: 93	Seleccionados: 1	Seleccionados: 1
Total: 93	Total: 2	Total: 22

#### ▶ Elija forma de presentación de la tabla

<b>Tipo de dato</b> ▾	◀
<b>Periodo</b> ▾	◀
<b>Clases</b> ▶	...
Decimales a mostrar: Por defecto	

#### ▶ Notas

Total: 93 series y 93 datos

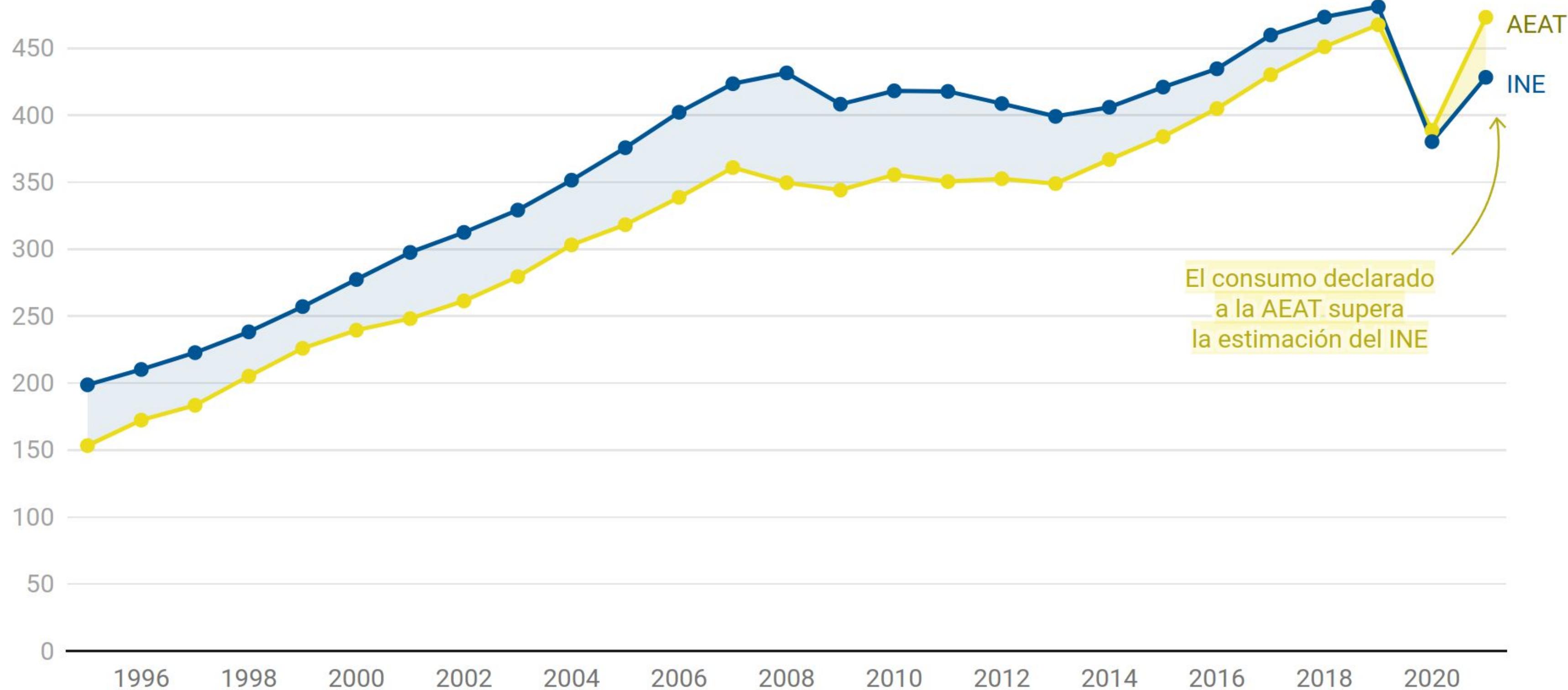
Consultar selección

Consultar todo



# El consumo declarado en 2021 supera al estimado por el INE

Evolución del consumo sujeto a IVA según la estimación del PIB (INE) y los datos de consumo declarado a la Agencia Tributaria.  
Datos en miles de millones de euros



El consumo declarado  
a la AEAT supera  
la estimación del INE

Muchos datos, fragmentados y desactualizados.

Para usarlos, tenemos que repetir procesamiento.

Juntar datos no es divertido ni fácil.

Pocas herramientas para colaborar y mejorar el ecosistema.

Con **datos** se pueden tomar  
mejores decisiones.

**Conectar datos entre sí** los  
hace más útiles.

# DATADEX

The Open Data Platform for your Community's Open Data.



Open-source, serverless, and local-first Data Platform to improve how communities coordinate.

[Get the Data](#)[Contribute](#)[Request Dataset](#)

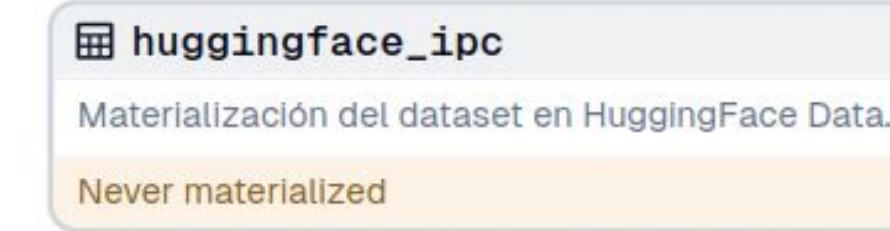
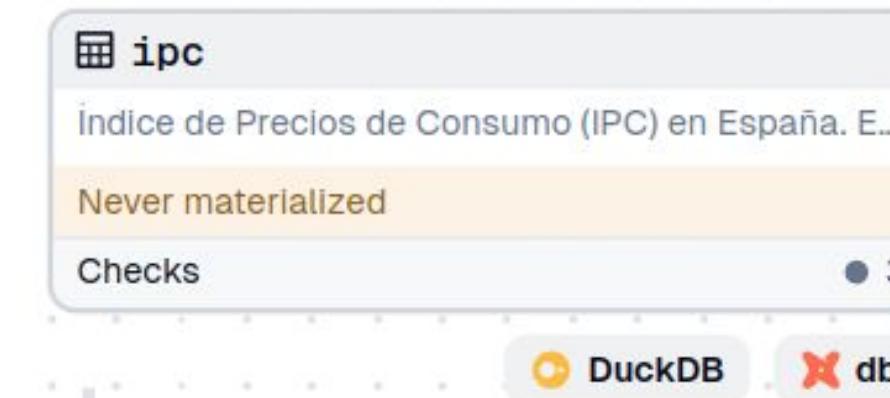
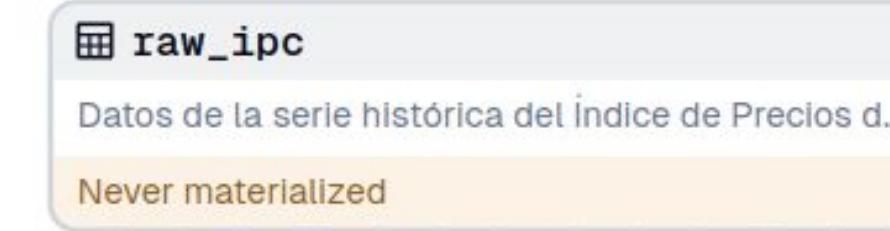
# Barefoot Data Portals

The Open Data Platform for your Community's Open Data.

Open-source, serverless, and local-first. Platform to improve how communities coordinate.

[Get the Data](#) [Contribute](#) [Request Dataset](#)



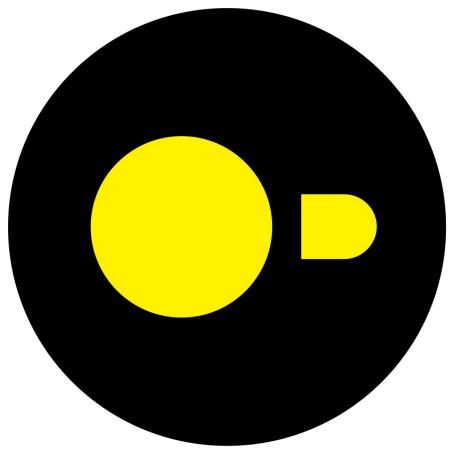




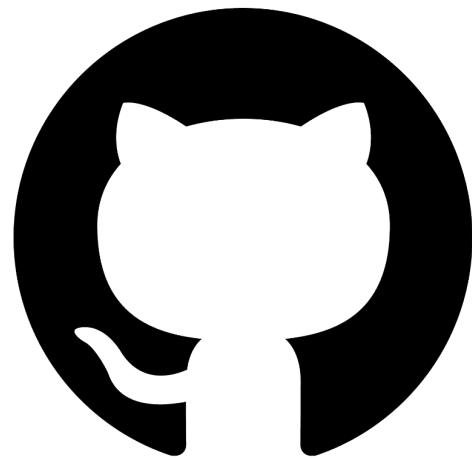
**dagster**



**dbt**

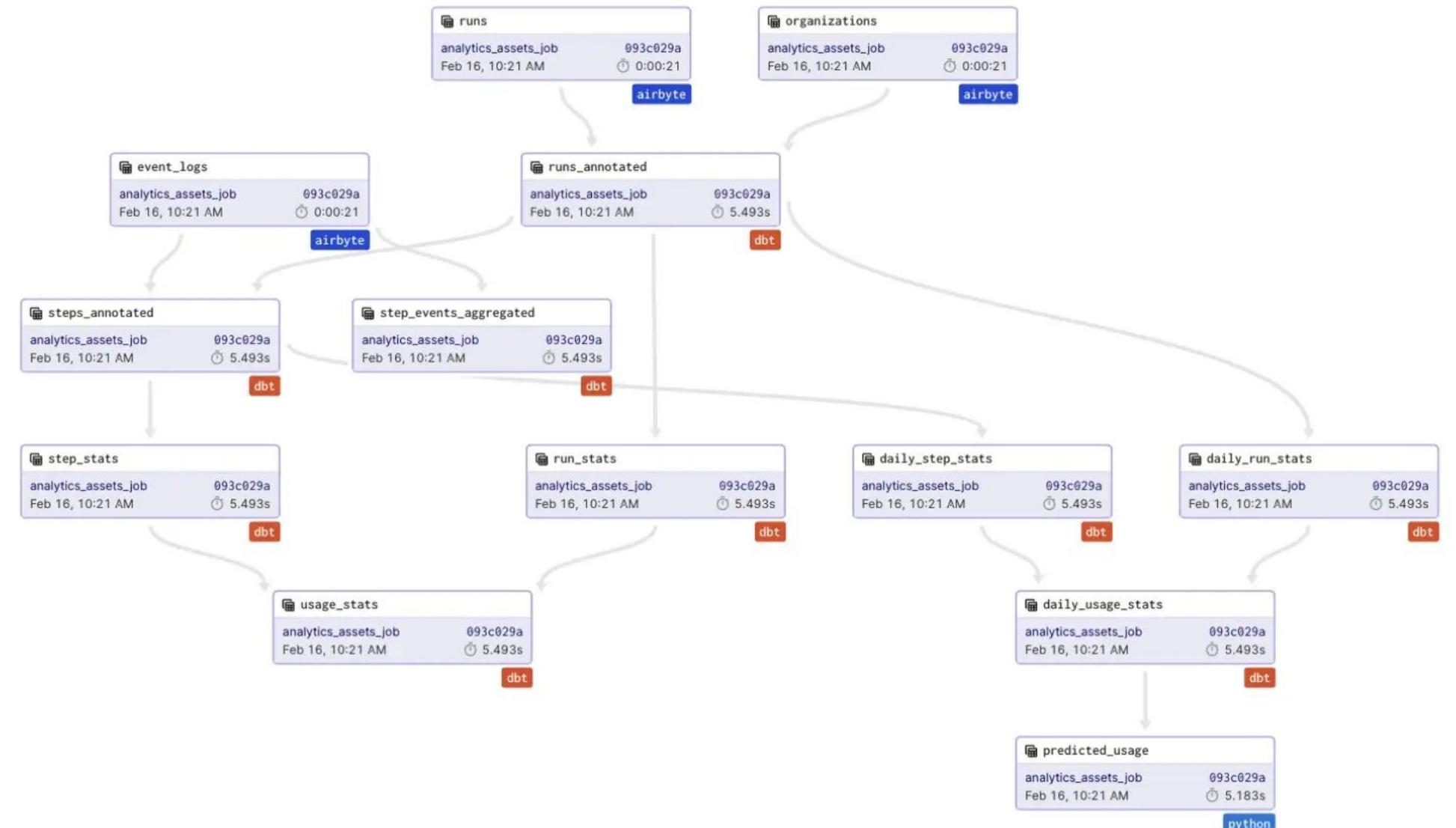
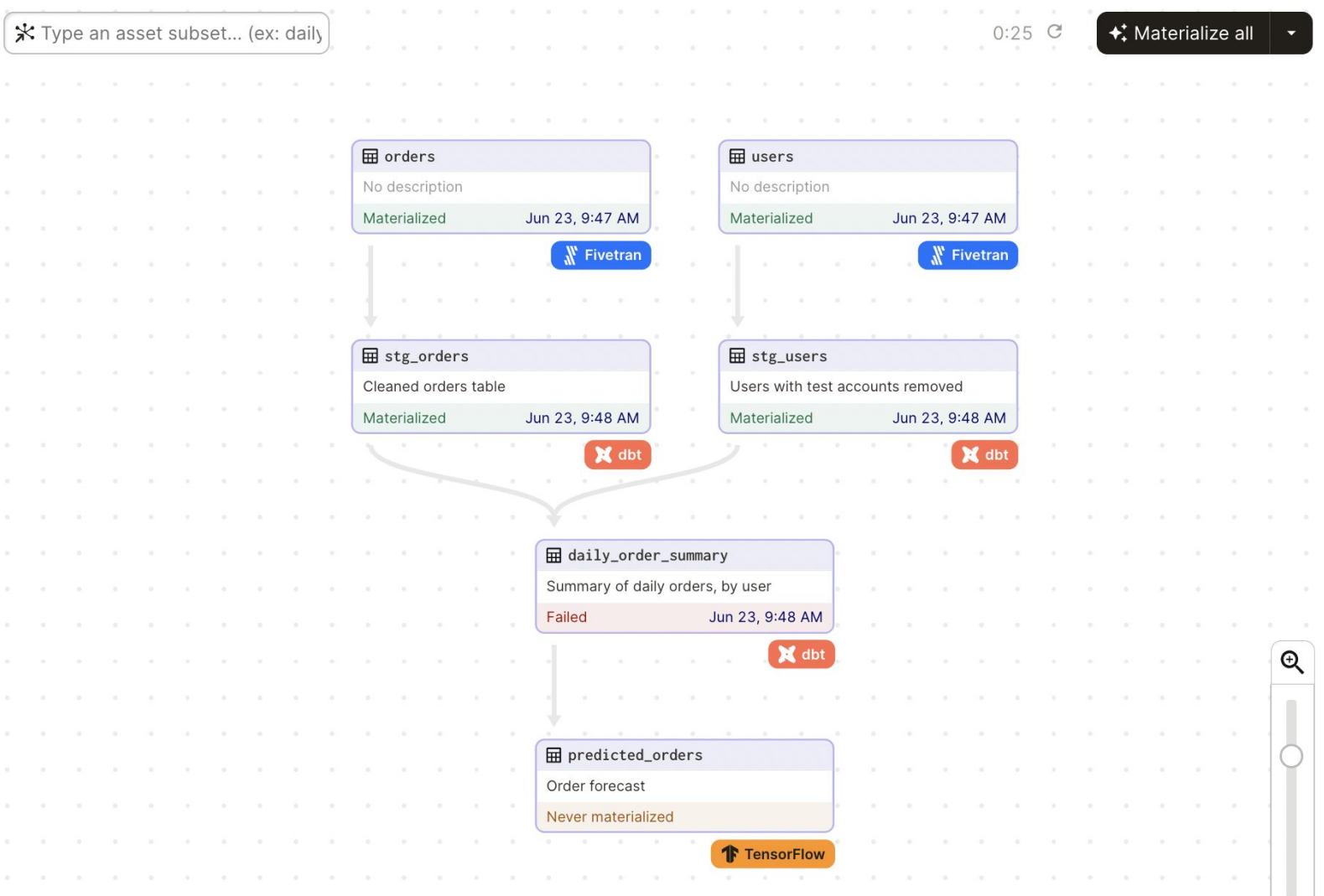
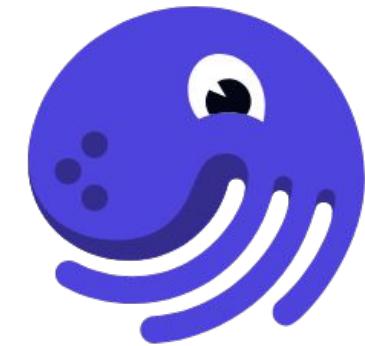


**DuckDB**



**GitHub**

# dagster





**dbt**

The screenshot shows the dbt Cloud interface for a project named "Fishtown Analytics". The left sidebar displays the project structure with a tree view of files and folders. The main area shows the SQL code for the "fct\_subscription\_transactions.sql" file. The code defines a materialized table with a source from "subscription\_transactions\_typed" and a window specification for calculating customer first month, month difference, and starting revenue. The right side of the interface shows the results of a recent run, including a table of data with columns ID, DATE\_MONTH, REVENUE, and REVENUE\_CHA, and a summary bar at the bottom indicating 100 Rows and 2.0 sec.

Project: Fishtown Analytics

branch: master

internal-analytics

- .github
- analytics
- data
- ddl
- macros

models

- export
- marts
- community
- consulting
- dbt

finance

- mrr
- utils

fct\_subscription\_transactions.sql

schema.yml

transactions

product

sales

staging

util

fct\_subscription\_transactions.sql

commit... commit...

compile run limit 100 save

1 {{ config(materialized = 'table') }}

2 with source as (

3 select \* from {{ref('subscription\_transactions\_typed')}}

4 ),

5 windows as (

6 select

7 \*,

8 min(date\_month) over (

9 partition by customer\_id

10 ) as customer\_first\_month,

11 datediff(month,

12 min(date\_month) over (partition by customer\_id),

13 date\_month

14 ) as customer\_month,

15 first\_value(revenue) over (

16 partition by customer\_id

17 order by date\_month

18 rows between unbounded preceding and unbounded following

19 ) as customer\_starting\_revenue

20 from source

21

22

23

24

25

26

27

28

29

30

31

32

ID	DATE_MONTH	REVENUE	REVENUE_CHA
58aa5f22ef7488e9a55349708002ae46	2019-01-01	150	150
cea7be635f2833a0aa59ece3ee1f8e08	2019-07-01	4000	4000
7089a5e0b5624543cccb27ffb10f96f5	2018-03-01	6000	3000
c23b85ff2b8f604642b4d630671e8251	2018-01-01	100	100
6423414d9ec902e63b6fd2bbc11ed538	2019-08-01	188.59	88.59
03f56e009103c2f7a0992940fcda57bbe	2019-07-01	2000	2000
03d8ce0c52b5e70dde4c47a0475d727d	2019-04-01	350.32	90.32
ab4ffc603ec4f6d69174707f2d66bbc8	2019-05-01	11000	5000
9e0432e3ab08493457984f09f1bd5a91	2017-06-01	6000	6000
4162c7d0a692e0a424392f14f61f8bf4	2018-07-01	5000	5000
947b01663b22f69aceeed92b9d616339	2019-07-01	100	100
894f6533f06ec532b1a49ad255ad0b3b	2018-07-01	1200	1200
3165067175c0fb38f9b959aa6be6bdd0	2018-04-01	4800	1800
e6d349c08869e455e0d3b1bb8d7a4910	2019-06-01	100	100
2a0536ff1283dc3d9953c254b45d97e6	2018-09-01	3000	3000
1ba92b8226c34e50b6bc790c4111300e	2018-11-01	110.03	10.03
0d466a073191eb7313cc9220cb6467ef02	2019-12-01	101.52	1.52

File Run • Success

100 Rows 2.0 sec

Runs dbt run ready

The screenshot shows the dbt Oracle interface with the following details:

- EXPLORER** pane on the left:
  - DBT-ORACLE project expanded.
  - Files listed include: countries.sql, eu\_direct\_sales\_channels\_promo\_costs.sql, direct\_sales\_channel\_promo\_cost.sql, exposures.yml, income\_levels.sql, internet\_sales\_channel\_customers.sql, people.sql, promotion\_costs.sql, sales\_cost.sql, sales\_internet\_channel.sql, schema.yml, union\_customer\_sales.sql, us\_product\_sales\_channel\_ranking\_append.sql, us\_product\_sales\_channel\_ranking.sql, us\_seed\_customers.sql.
  - Seeds folder expanded, containing seeds and snapshots.
  - promotion\_costs.sql file is selected.
  - target, compiled, and run folders.
  - OUTLINE section shows: No symbols found in document 'us\_product\_sales\_channel\_ranking.sql'.
- Code Editor** pane in the center:

```
! profiles.yml M  internet_sales_channel_customers.sql  ●  us_product_sales_channel_ranking.sql X  dbt-oracle
dbt_adbs_test_project > models > us_product_sales_channel_ranking.sql
14     limitations under the License.
15   #}
16   {
17     config(
18       materialized='incremental',
19       unique_key='group_id',
20       parallel=4,
21       table_compression_clause='COLUMN STORE COMPRESS FOR QUERY LOW')
22   }
23
24   SELECT prod_name, channel_desc, calendar_month_desc,
25     {{ snapshot_hash_arguments(['prod_name', 'channel_desc', 'calendar_month_desc']) }} AS group_id,
26     TO_CHAR(SUM(amount_sold), '9,999,999,999') SALES$,
27     RANK() OVER (ORDER BY SUM(amount_sold)) AS default_rank,
28     RANK() OVER (ORDER BY SUM(amount_sold) DESC NULLS LAST) AS custom_rank
29   FROM {{ source('sh_database', 'sales') }}, {{ source('sh_database', 'products') }}, {{ source('sh_database', 'times') }}, {{ source('sh_database', 'channels') }}, {{ source('sh_database', 'customers') }}
30   WHERE sales.prod_id=products.prod_id AND sales.cust_id=customers.cust_id
31     AND customers.country_id = countries.country_id AND sales.channel_id=channels.channel_id
32     AND country_iso_code='US'
33
34   {% if is_incremental() %}
35
36     AND times.calendar_month_desc > (SELECT MAX(calendar_month_desc) FROM {{ this }})
37
38   {% endif %}
39
40   GROUP BY prod_name, channel_desc, calendar_month_desc
41
42
43
44
```
- TIMELINE** pane at the bottom left.
- Bottom status bar: Ln 1 Col 1 Spaces: 2 UTF-8 LF (à SQL)

DuckDB Shell

```

Elapsed: 832 ms
Recent downloads
customer.parquet
1,167 KB - 1 minute ago
Show all downloads
duckdb> PRAGMA tpch(7);

```

supp_nation	cust_nation	l_year	revenue
FRANCE	GERMANY	1995	4637235.1501
FRANCE	GERMANY	1996	5224779.5736
GERMANY	FRANCE	1995	6232818.7037
GERMANY	FRANCE	1996	5557312.1121

```

Elapsed: 57 ms
duckdb> COPY customer TO 'customer.parquet';

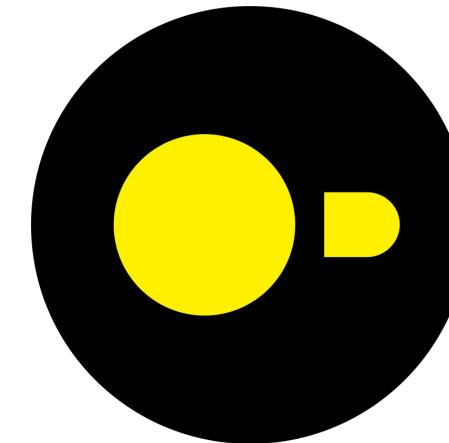
```

Count
15000

```

Elapsed: 139 ms
duckdb> .files download customer.parquet
Copied file: customer.parquet
duckdb>

```



# DuckDB

```

input_df = pd.DataFrame.from_dict({"i":[1, 2, 3],
                                    "j":["Apple", "Mango", "Oranges"]})

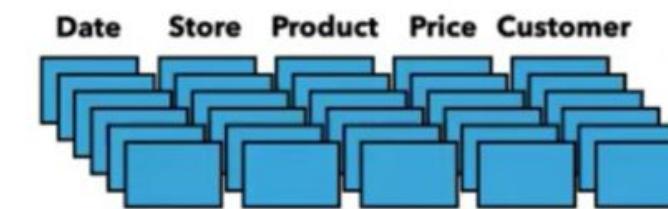
%sql output_df << SELECT sum(i) as total_i FROM input_df
output_df

```

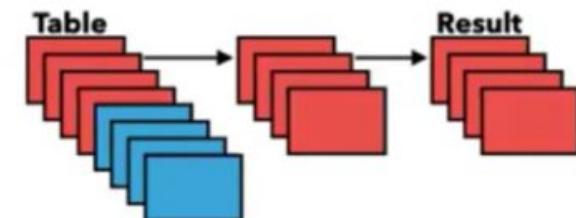
total_i	
0 6.0	

## CWI DuckDB at a Glance

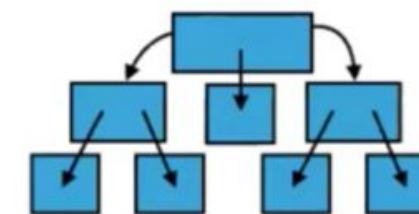
### Column-Store



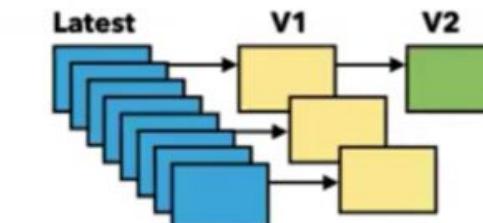
### Vectorized Processing



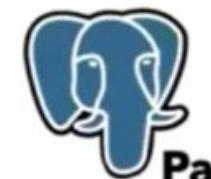
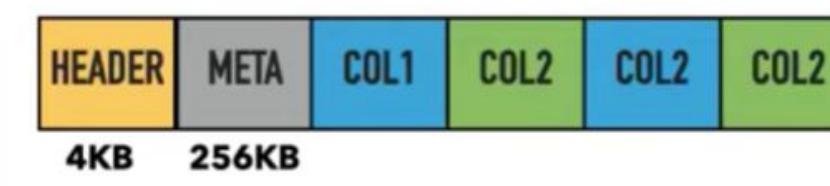
### ART Index



### Multi-Version Concurrency Control



### Single-File Storage



Parser

## Open a pull request

Create a new pull request by comparing changes across two branches. If you need to, you can also [compare across forks](#).

base: main ▾ ← compare: my-patch-1 ✓ Able to merge. These branches can be automatically merged.

Choose a head ref

Update CONT  
my

Write Prev  
Branches Tags  
✓ my-patch-1

Leave a comment  
myarb-patch-1

my

jjallaire / jjallaire.github.io Public Pin Unwatch 1 Fork 0 Star 1

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

General Access Collaborators Moderation options

GitHub Pages

GitHub Pages is designed to host your personal, organization, or project pages from a GitHub repository.

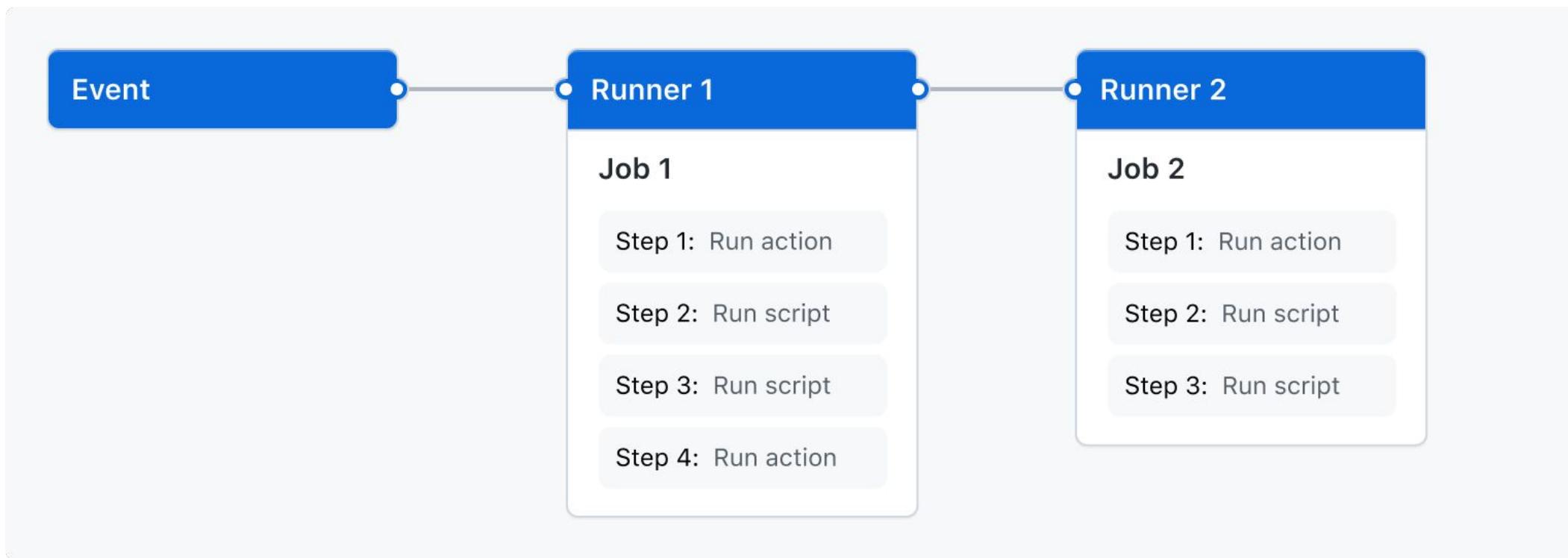
Your site is published at <https://jjallaire.github.io/>

Source Branch: gh-pages / (root) Save

Code and automation Branches Tags Actions Webhooks Environments

Custom domain Custom domains allow you to serve your site from a domain other than jjallaire.github.io. Learn more.

Pages Save Remove

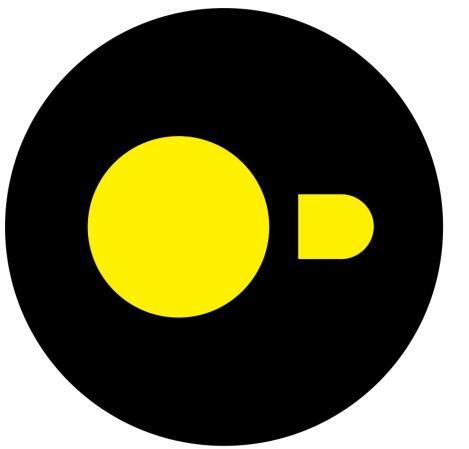




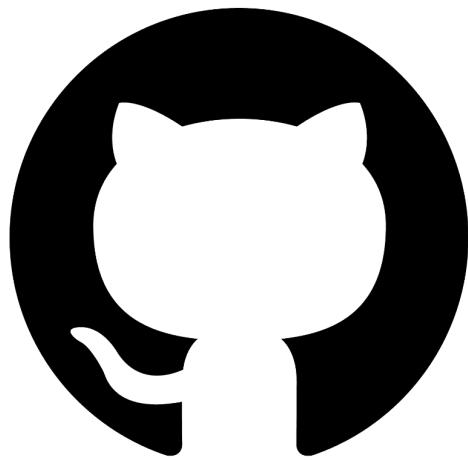
**dagster**



**dbt**



**DuckDB**



**GitHub**



Apache  
Airflow





# datadex

(Public)

[Edit Pins](#)[Unwatch 4](#)[Fork 15](#)[Star 246](#)[main](#)[1 Branch](#)[9 Tags](#)[Go to file](#)

t

[Add file](#)[Code](#)

davidgasquez feat: ✨ Add multi-commit support in DatasetPublisher



93da061 · 2 months ago

329 Commits

[.devcontainer](#)

chore: 🛡️ remove slqfmt [skip ci]

6 months ago

[.github/workflows](#)

chore: 🔧 adjust concurrency and python version in ci

5 months ago

[data](#)

fix: 🐛 don't ignore the data folder

9 months ago

[datadex](#)

feat: ✨ Add multi-commit support in DatasetPublisher

2 months ago

[dbt](#)

feat: ⚡ enable spain weather assets

5 months ago

[portal](#)

feat: 🚀 Improve performance and UX with multiple enh...

5 months ago

[.editorconfig](#)

chore: 🔧

last year

[.gitignore](#)

refactor:💡 consolidate indicators assets under polars a...

5 months ago

[Dockerfile](#)

feat: 🐦 Upgrade Python and refactor for polars integrati...

5 months ago

[LICENSE](#)

Initial commit

2 years ago

[Makefile](#)

feat: 🐛 update code to use new dbt classes

2 months ago

[README.md](#)

docs: 📝 add lung-sarg [skip ci]

2 months ago

[pyproject.toml](#)

feat: ✨ Add multi-commit support in DatasetPublisher

2 months ago

[README](#)

MIT license



# DATADEX

## About

Serverless and local-first Open Data Platform

[datadex.datonic.io](#)

sql open-data dbt quarto  
duckdb

Readme

MIT license

Activity

Custom properties

246 stars

4 watching

15 forks

## Contributors 5



## Languages

Jupyter Notebook 75.2%

Python 23.2% Other 1.6%

```

99     @asset()
100    def spain_aemet_stations_data(aemet_api: AEMETAPI) -> pl.DataFrame:
101        """
102            Spain AEMET stations data.
103        """
104
105        df = pl.DataFrame(aemet_api.get_all_stations())
106        df.with_columns(pl.col("indsinop").cast(pl.Int32, strict=False).alias("indsinop"))
107
108        # Clean latitud and longitud
109        def convert_to_decimal(coord):
110            degrees = int(coord[:-1][:2])
111            minutes = int(coord[:-1][2:4])
112            seconds = int(coord[:-1][4:])
113            decimal = degrees + minutes / 60 + seconds / 3600
114            if coord[-1] in ["S", "W"]:
115                decimal = -decimal
116            return decimal
117
118        df = df.with_columns(
119            [
120                pl.col("latitud").map_elements(convert_to_decimal).alias("latitud"),
121                pl.col("longitud").map_elements(convert_to_decimal).alias("longitud"),
122            ]
123        )
124
128        @asset()
129        def spain_aemet_weather_data(
130            context: AssetExecutionContext, aemet_api: AEMETAPI
131        ) -> pl.DataFrame:
132            """
133            Spain weather data since 1940.
134            """
135
136            start_date = datetime(1940, 1, 1)
137            end_date = datetime.now()
138
139            r = aemet_api.get_weather_data(start_date, end_date)
140
141            df = pl.DataFrame()
142            for d in r:
143                ndf = pl.DataFrame(d)
144                df = pl.concat([df, ndf], how="diagonal_relaxed")
145
146            df = df.with_columns(pl.col("fecha").str.strptime(pl.Date, format="%Y-%m-%d"))
147
148            float_columns = [
149                "prec",
150                "presMax",
151                "presMin",
152                "racha",
153                "sol",
154                "tmax",
155                "tmed",
156                "tmin",
157                "velmedia",
158            ]
159
160            df = df.with_columns(
161            [
162                pl.col(col).str.replace(",", ".").cast(pl.Float64, strict=False)
163                for col in float_columns
164            ]
165        )

```

# ESTACIONES TIEMPO

```
128     @asset()
129     def spain_aemet_weather_data(
130         start_date: pl.Date,
131         end_date: pl.Date,
132         aemet_api: AEMETAPI
133     ) -> pl.DataFrame:
134
135     return (
136         select
137             cast(w.fecha as date) as fecha,
138             w.indicativo,
139             w.nombre,
140             w.provincia,
141             s.latitud,
142             s.longitud,
143             cast(w.altitud as int) as altitud,
144             cast(w.tmed as float) as tmed,
145             cast(w.prec as float) as prec,
146             cast(w.tmin as float) as tmin,
147             w.horatmin,
148             cast(w.tmax as float) as tmax,
149             w.horatmax,
150             cast(w.dir as int) as dir,
151             cast(w.velmedia as float) as velmedia,
152             cast(w.racha as float) as racha,
153             w.horaracha,
154             cast(w.presMax as float) as presMax,
155             w.horaPresMax,
156             cast(w.presMin as float) as presMin,
157             w.horaPresMin,
158             cast(w.hrMedia as int) as hrMedia,
159             cast(w.hrMax as int) as hrMax,
160             w.horaHrMax,
161             cast(w.hrMin as int) as hrMin,
162             w.horaHrMin,
163             cast(w.sol as float) as sol
164         from {{ source('main', 'spain_aemet_weather_data') }} as w
165         left join {{ source('main', 'spain_aemet_stations_data') }} as s
166             on w.indicativo = s.indicativo
167         order by w.fecha, w.indicativo
168         .cast(pl.Float64, strict=False)
169     )
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
449
450
451
452
453
454
455
456
457
458
459
459
460
461
462
463
464
465
466
467
468
469
469
470
471
472
473
474
475
476
477
478
479
479
480
481
482
483
484
485
486
487
488
489
489
490
491
492
493
494
495
496
497
498
499
499
500
501
502
503
504
505
506
507
508
509
509
510
511
512
513
514
515
516
517
518
519
519
520
521
522
523
524
525
526
527
528
529
529
530
531
532
533
534
535
536
537
538
539
539
540
541
542
543
544
545
546
547
548
549
549
550
551
552
553
554
555
556
557
558
559
559
560
561
562
563
564
565
566
567
568
569
569
570
571
572
573
574
575
576
577
578
579
579
580
581
582
583
584
585
586
587
588
589
589
590
591
592
593
594
595
596
597
598
599
599
600
601
602
603
604
605
606
607
608
609
609
610
611
612
613
614
615
616
617
618
619
619
620
621
622
623
624
625
626
627
628
629
629
630
631
632
633
634
635
636
637
638
639
639
640
641
642
643
644
645
646
647
648
649
649
650
651
652
653
654
655
656
657
658
659
659
660
661
662
663
664
665
666
667
668
669
669
670
671
672
673
674
675
676
677
678
679
679
680
681
682
683
684
685
686
687
688
689
689
690
691
692
693
694
695
696
697
698
699
699
700
701
702
703
704
705
706
707
708
709
709
710
711
712
713
714
715
716
717
718
719
719
720
721
722
723
724
725
726
727
728
729
729
730
731
732
733
734
735
736
737
738
739
739
740
741
742
743
744
745
746
747
748
749
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
769
770
771
772
773
774
775
776
777
778
779
779
780
781
782
783
784
785
786
787
788
789
789
790
791
792
793
794
795
796
797
798
799
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
848
849
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
878
879
879
880
881
882
883
884
885
886
887
888
889
889
890
891
892
893
894
895
896
897
898
899
899
900
901
902
903
904
905
906
907
908
909
909
910
911
912
913
914
915
916
917
918
919
919
920
921
922
923
924
925
926
927
928
929
929
930
931
932
933
934
935
936
937
938
939
939
940
941
942
943
944
945
946
947
948
949
949
950
951
952
953
954
955
956
957
958
959
959
960
961
962
963
964
965
966
967
968
969
969
970
971
972
973
974
975
976
977
978
979
979
980
981
982
983
984
985
986
987
987
988
989
989
990
991
992
993
994
995
996
997
998
999
999
1000
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1079
1080
1081
1082
1083
1084
1085
1086
1087
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1097
1098
1099
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1187
1188
1189
1189
1190
1191
1192
1193
1194
1195
1196
1196
1197
1198
1199
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1287
1288
1289
1289
1290
1291
1292
1293
1294
1295
1296
1297
1297
1298
1299
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1387
1388
1389
1389
1390
1391
1392
1393
1394
1395
1396
1396
1397
1398
1399
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1487
1488
1489
1489
1490
1491
1492
1493
1494
1495
1496
1496
1497
1498
1499
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1587
1588
1589
1589
1590
1591
1592
1593
1594
1595
1596
1596
1597
1598
1599
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1678
1
```

# TEMPO.

# ESTACIONES



Split (1)

train · 8.11M rows

Ordered by fecha (desc) ×

Share results 8,114,317 rows

fecha	indicativo	nombre	provincia	latitud	longitud	altitud	tmed	prec	tmin	horatmin	tmax	horatmax
	unknown	large_string · classes	large_string · classes	large_string · classes	float64	float64	int32	float32	float32	large_string · lengths	float32	large_string · lengths
		947 values	934 values	53 values	27.7	43.8	1	0	-30	5	-50	5
"2024-09-15"	0016A	REUS AEROPUERTO	TARRAGONA	41.145	1.163611	71	19	0	10	05:35	28	12:38
"2024-09-15"	0034X	VALLS	TARRAGONA	41.293056	1.260833	233	18.4	0	8.8	05:14	28	13:58
"2024-09-15"	0042Y	TARRAGONA	TARRAGONA	41.123889	1.249167	55	19.799999	0	11.7	06:10	28	12:07
"2024-09-15"	0061X	PONTONS	BARCELONA	41.416944	1.519167	632	16.5	0	7.4	02:24	25.6	14:43
"2024-09-15"	0066X	VILAFRANCA DEL PENEDÈS	BARCELONA	41.330278	1.676944	177	19.4	0	10.3	05:31	28.5	14:37
"2024-09-15"	0073X	SITGES	BARCELONA	41.243889	1.8525	58	19.6	0	12.6	04:38	26.6	13:28
"2024-09-15"	0076	BARCELONA AEROPUERTO	BARCELONA	41.292778	2.07	4	20.200001	0	14	05:30	26.299999	14:13
"2024-09-15"	0092X	BERGA	BARCELONA	42.101389	1.8575	682	16.299999	0	7.6	05:06	25	14:46
"2024-09-15"	0106X	BALSARENY	BARCELONA	41.866389	1.8725	361	16	0	4.7	05:14	27.299999	14:14
"2024-09-15"	0114X	PRATS DE LLUÇÀNÈS	BARCELONA	42.006944	2.026667	700	15.8	0	6	05:54	25.5	15:31
"2024-09-15"	0120X	MOIÀ	BARCELONA	41.813333	2.095278	742	17	0	10.3	05:23	23.700001	13:55
"2024-09-15"	0149X	MANRESA	BARCELONA	41.72	1.840278	291	17.4	0	7.4	06:10	27.5	14:40
"2024-09-15"	0158X	MONISTROL DE MONTserrat	BARCELONA	41.594444	1.839167	738	17	0	11.7	05:23	22.299999	13:24
"2024-												

main

1 Branch

0 Tags

Go to file

t

Add file

Code

## About

No description, website, or topics provided.

Activity

5 stars

1 watching

1 fork

[Report repository](#)

## Releases

No releases published

## Sponsor this project



simonw Simon Willison

Sponsor

[Learn more about GitHub Sponsors](#)

## Packages

No packages published

Automated Sat Oct 12 16:04:47 UTC 2024

1187e25 · 24 minutes ago 177 Commits

.github/workflows

Scrape outage.tecoenergy.com

2 days ago

duke-energy.app

Sat Oct 12 16:04:47 UTC 2024

24 minutes ago

fplmaps.com

Sat Oct 12 16:04:47 UTC 2024

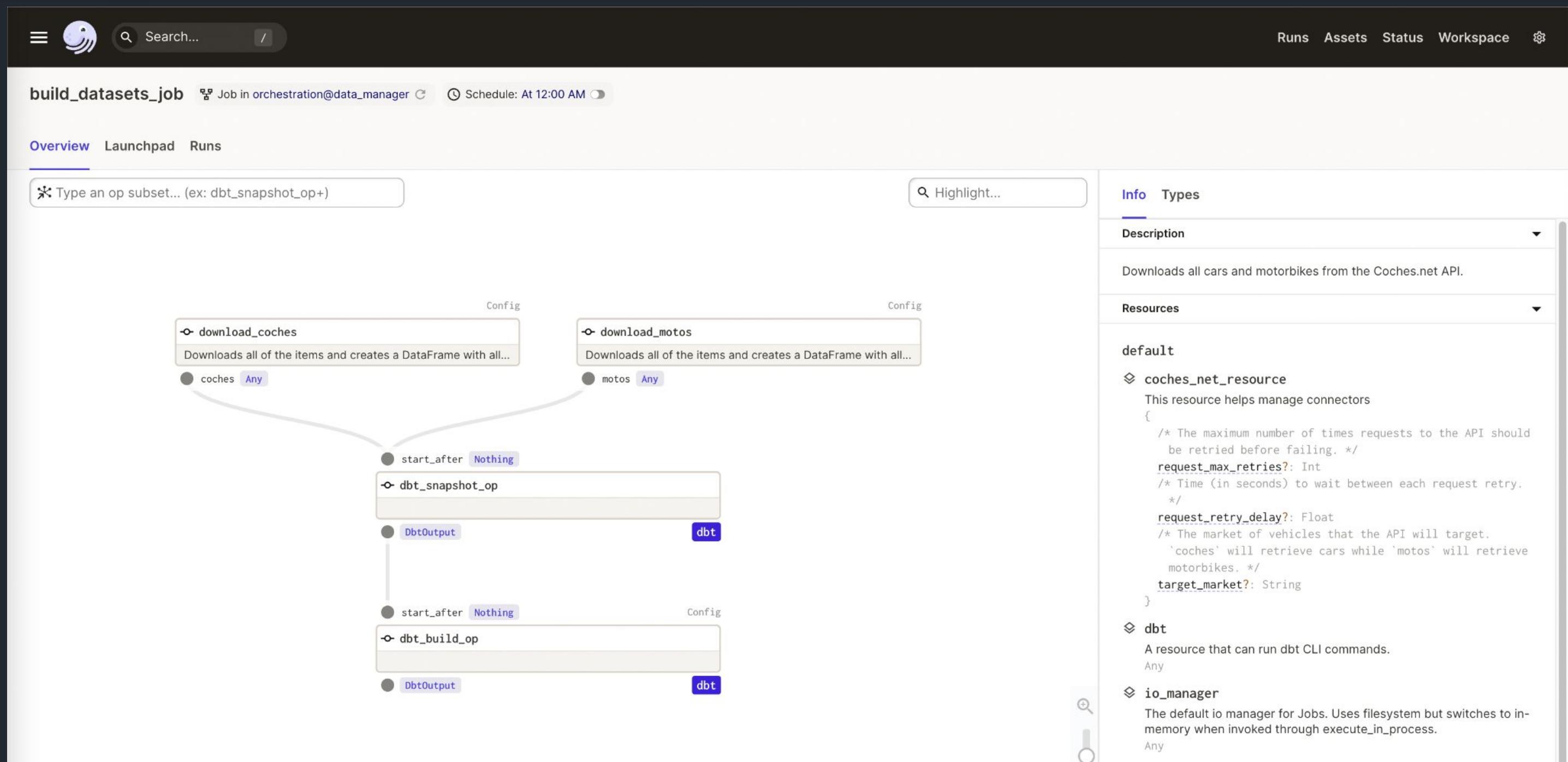
24 minutes ago

outage.tecoenergy.com

Sat Oct 12 16:04:47 UTC 2024

24 minutes ago

```
12   jobs:
13     fetch:
14       runs-on: ubuntu-latest
15       steps:
16         - uses: actions/checkout@v4
17         - name: Scrape fplmaps.com
18           run: |
19             mkdir -p fplmaps.com
20             if ! curl 'https://www.fplmaps.com/customer/outage/StormFeedRestoration.json' \
21               -H 'User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10.15; rv:131.0) Gecko/20100101 Firefox/131.0' \
22               -H 'Accept: application/json, text/javascript, */*; q=0.01' \
23               | jq > fplmaps.com/StormFeedRestoration.json; then
24               echo "Failed to scrape fplmaps.com, but continuing..."
25             fi
```



# coches-net-dashboard

Coches.net dashboard is an application made by [@franloza](#) and [@Jorjatorz](#) to visualize Spanish car and motorcycle market through the data obtained from the websites [coches.net](#) and [motos.coches.net](#).

This project is composed of a data pipeline run with `Dagster`, which extract data from `coches.net`, apply some transformations using `dbt`, and stores the data in a DuckDB database. The data can be analyzed with a dashboard built with `Dash`.

This application should not be considered suitable for production, and it's intended to be used only locally for analytical purposes.

## Contributors 2

[Jorjatorz](#) Jorge Sanchez

[franloza](#) Fran Lozano

## Languages

Python 97.6% Dockerfile 2.4%

# The Public Utility Data Liberation (PUDL) Project

[ACCESS PUDL DATA](#)[PUDL ON GITHUB](#)

Electric utilities report a huge amount of information to the US government and other public agencies. This includes yearly, monthly, and even hourly data about fuel burned, electricity generated, operating expenses, power plant usage patterns and emissions. Unfortunately, much of this data is not released in well documented, ready-to-use, machine readable formats. Data from different agencies tends not to be standardized or easily used in tandem. Several commercial data services clean, package, and re-sell this data, but at prices which are too high to be accessible to many smaller stakeholders.

The [Public Utility Data Liberation \(PUDL\) project](#) takes information that's already publicly *available*, and makes it publicly *usable*, by cleaning, standardizing, and cross-linking utility data from different sources in a single database. Thus far our primary focus has been on fuel use, generation, operating costs, and operation history. As of November, 2022 PUDL integrates data from:

- [EIA Form 860](#): 2001-2023
- [EIA Form 860m](#): 2023-06
- [EIA Form 861](#): 2001-2022
- [EIA Form 923](#): 2001-2022
- [EPA Continuous Emissions Monitoring System \(CEMS\)](#): 1995-2022
- [FERC Form 1](#): 1994-2022
- [FERC Form 714](#): 2006-2020

# Research and data to make progress against the world's largest problems.

Try "Life expectancy", "Poverty Nigeria Vietnam", "CO2 France"...



13,189 charts across 119 topics — All free: open access and open source

## OUR MISSION

### What do we need to know to make the world a better place?

To make progress against the pressing problems the world faces, we need to be informed by the best research and data.

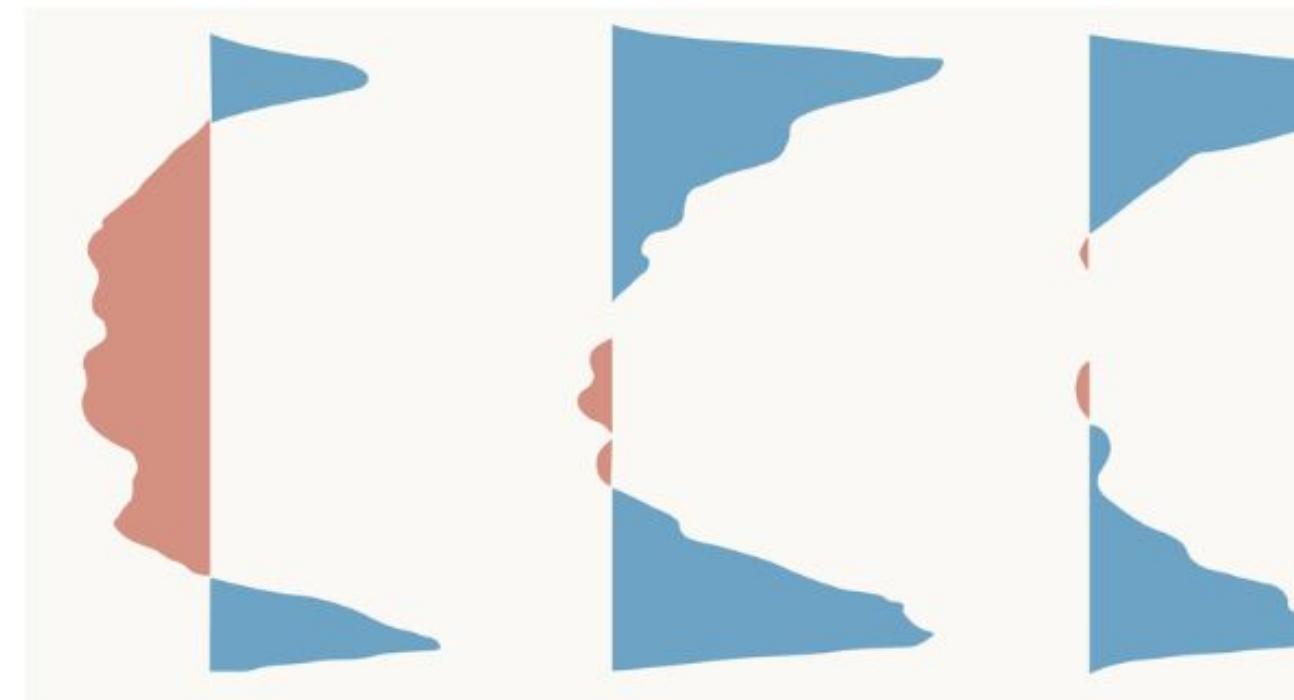
Our World in Data makes this knowledge accessible and understandable, to empower those working to build a better world.

[Read about our mission →](#)

[Subscribe to our newsletter](#)

We are a non-profit — all our work is free to use and open source. Consider supporting us if you find our work valuable.

## FEATURED WORK



ARTICLE · 15 MIN READ

### How will climate change affect crop yields in the future?

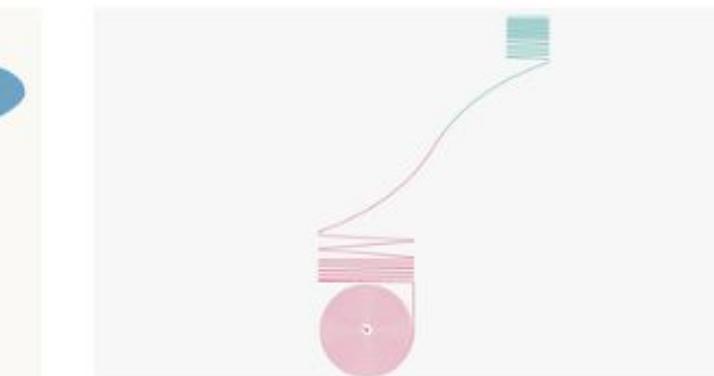
Maize yields could see significant declines, but wheat could increase. Impacts across the world will be very different.

Hannah Ritchie

#### REVISED & UPDATED

#### Cancer

Saloni Dattani, Veronika Samborska, Hannah

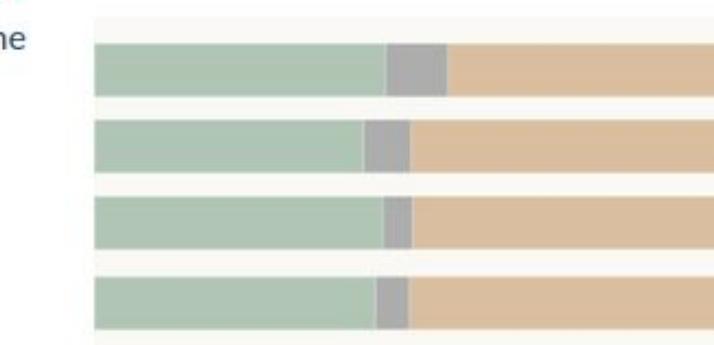


#### FROM OUR CLASSICS

### Technology over the long run: zoom out to see how dramatically the world can change within a lifetime

It is easy to underestimate the magnitude of this change. Understanding this can help us see how different the world could be in the future.

Max Roser



#### ARTICLE · 12 MIN READ

### Crop yields have increased dramatically in recent decades, but

[Design principles](#)[Our journey](#)[Fundamentals](#)[Computational graph](#)[ETL model](#)[Data model](#)[Features and constraints](#)[The DAG](#)[The URI](#)[ETL steps](#) [Other steps](#)

# ETL model

[On this page](#)[Snapshots](#)[Datasets](#)[Grapher views](#)

This section presents a high-level view of our ETL processes.

In general, an ETL consists of three main processes:

- **Extract:** Downloads data from a source or multiple sources.
- **Transform:** Processes and combines the downloaded data so it is usable.
- **Load:** Moves the transformed data to a data repository (e.g. database).

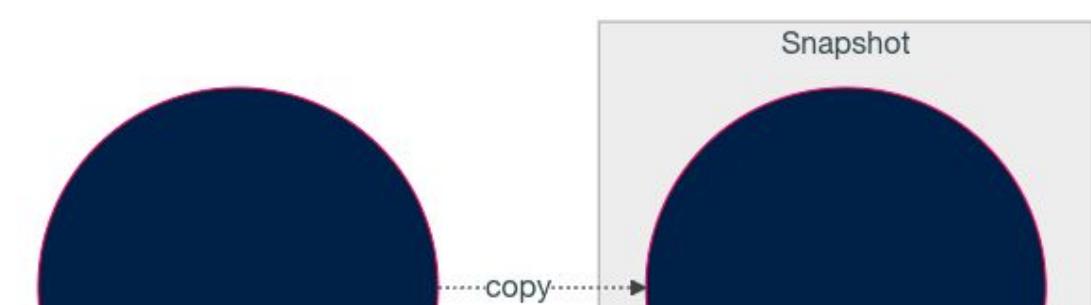
Our ETL follows these same principles, with different names, and has the following three main building blocks:

- **Snapshots:** We download files from external sources and store them as snapshots.
- **Datasets:** Next, we curate the downloaded data to have clean datasets. These datasets can depend on snapshots, but also on other curated datasets.
- **Grapher views:** We adapt the curated datasets for Grapher, which are then loaded to Grapher.

As a consequence, nodes in [our DAG](#) can come in different flavours.

## Snapshots

Snapshots are edge nodes in the computational graph. They represent a copy of a upstream data file — preserving its original format — at a particular point in time. That is, they are entry points to the ETL, and therefore don't have dependencies within the computational graph.



Plataformas de Datos que **no**  
**añaden datos nuevos**, sino que  
juntan, limpian datos existentes,  
haciéndolos **útiles** y **fáciles** de  
acceder para una comunidad.

1. Un repositorio donde poner el código para colaborar y un sitio para ejecutarlo.
2. Trabajar con estándares y tecnologías modernas (ahorra tiempo y dinero).
3. Publicar ficheros estáticos (no APIs) que se actualizan automáticamente.
4. Trabajar cerca de la comunidad para saber cuáles son sus problemas y donde pueden ayudar los datos.

# ¡Gracias!

¿Preguntas? ¿Feedback?



@davidgasquez

# DATALIA

*Datos, sin complicaciones.* 

---

Datalia es una plataforma de datos abiertos a nivel de España con el objetivo de unificar y armonizar información proveniente de diferentes fuentes.

## Datasets

Gracias a la colaboración de la comunidad, Datalia ofrece una variedad de conjuntos de datos. A continuación, se muestran algunos de los conjuntos de datos disponibles.

- [IPC](#)

## Recursos

Algunos recursos que pueden ser de utilidad relacionados con datos abiertos y transparencia a nivel de España.

## Portales

- [INE.](#)
- [Datos Abiertos de España.](#)
- [DGT.](#) Tambien en [mapa](#).
- [AEMET.](#)
- [WikiData.](#)
- [ESRI.](#)
- [Datadista.](#)
- [Newtral.](#)
- [Spanish Origin Destination Data](#)
- [DataMarket.](#)