

OPEN DATA TUTORIAL

Opening and promoting use of budget data

Identify the source and format

Once the data to be published has been prioritized according to its potential users and their information demands, it is vital to locate the origin and format of the information.

Why does the origin and format of the data matter? Although a recent emphasis has been made in data being open by default, most of the times **the information systems containing the data to be opened were not created having its publication and/or reuse in mind, but to display the daily operations of the government instead.** Thus, working on its openness will involve analysis, debugging and creation of new processes to ensure its sustainability.

Therefore it is worth answering the following questions:

- Is the data in **printed format only**?
- Is all the data to be published in the **same place**?
- How many people are responsible for **generating, concentrating** and/or **publishing** the data?
- In what **format** is it: PDF, Word, Excel, CSV or another?
- Does the data come from an **information system**?
- If so, how can we **access** it?
- Is it all in the same system or is it necessary **to join data from various technological tools**?



Considering these questions helps establish times for the opening process; identify the areas involved; analyze the quality of the information; as well as defining the challenges and opportunities ahead.

TIP

Make an operative calendar in a collaborative fashion with the areas involved, to define specific dates and actions towards the publication of the information.

Some relevant considerations to identify the source and format of the data are:

1. Identify the actors involved

Spending information may be scattered amongst areas and even government agencies. Nonetheless, it is important to take into account that the more dispersed the information is and the more people are responsible for it, the more will the process of openness depend on the existing cooperation capacities.

Conversely, a set of data that is in a single place with less number of co-responsible areas, can become the spearhead of the opening process, especially if you have already identified public interest on it.

Likewise, consider that the areas responsible for the information may not be familiar with the concepts of data openness and at the same time, the people in charge of the openness may not be very familiar with the contents and structure of the data sets that are being published; so, it is essential to start a continuous two-way learning process, since the usability and quality of the resulting data set and its data dictionary rely on it.

2. Data format

In order to allow the information to be used effectively, it is necessary to consider how malleable or rigid the data is, whether it is ready for processing or if it can be consumed in whole and/or in parts. This will depend, at first glance, on whether the data is unstructured or structured, and from there, you will be able to determine the possibilities of use.

Unstructured data refers to disorganized information, that is, a conglomerate of information from which specific attributes cannot be extracted in an agile and/or automated way.

These data may be public or not, but it is characterized by its storing, for example: information obtained directly from beneficiaries of a service or product derived from official interactions between public servants is often in printed format; information digitized from printed or photographic material is stored as digital images (PDF, JPG, PNG); and reports can be found as text editor files or slide presentations (DOCX, PPTX), among others.

Because of our daily interaction with this formats, we know that it is not possible to obtain sums or any mathematical operation from the information as it is more complicated to process it, so, it will be necessary to transform the data and fit it into a structured format.

In contrast, **structured data** refers to information readable by machines (via open sourced or licensed software), since it has an easily recognizable arrangement. For example, if the structure of the information consists of rows and columns, the spreadsheet or linear display formats allow its agile processing (XLSX, CSV, etc.); while for tree structured-information, the most common formats for information hierarchy are JSON or XML.

If the data that will be published is in a structured format, the openness process can be speeded up.



SHCP

1.2 Producción y empleo

1.2.1 Producción

Durante el segundo trimestre de 2017, el PIB registró un crecimiento anual de 1.8 por ciento real. Al excluir el efecto estacional, el PIB creció a una tasa de 3.0 por ciento anual, y de 0.6 por ciento en términos trimestrales. A su interior, se observaron los siguientes resultados:

- La producción agropecuaria se incrementó a una tasa anual de 0.7 por ciento. Cifras ajustadas por estacionalidad indican que este sector tuvo una reducción trimestral de 1.9 por ciento.
- La actividad industrial tuvo una disminución anual de 1.1 por ciento. Dicho resultado fue consecuencia, principalmente, de la menor producción de la minería petrolera y las obras de ingeniería civil. Por su parte, la producción de manufacturas registró un aumento anual de 2.0 por ciento, debido a la expansión de los subsectores de equipos de transporte, equipo de computación, alimentos y fabricación de maquinaria y equipo, entre otros. Al eliminar el efecto estacional, la producción industrial no presentó cambio alguno respecto al trimestre previo.
- El sector servicios registró un crecimiento anual de 3.2 por ciento como resultado del desempeño de las actividades de servicios financieros, comercio, información en medios masivos, inmobiliarias y de alquiler, y transportes, correos y almacenamiento, principalmente. Cifras ajustadas por estacionalidad indican que los servicios aumentaron a una tasa trimestral de 0.8 por ciento.

PRODUCTO INTERNO BRUTO, 2015-2017 A*/
(Variación % anual)

	Enero-junio			2015			2016			2017		
	2015	2016	2017	I	II	III	I	II	III	I	II	III
Total	2.7	2.4	2.3	2.5	2.8	2.5	2.2	2.6	3.0	2.3	2.8	1.8
Agropecuaria	2.7	1.8	3.1	0.5	0.5	0.4	-0.6	1.9	5.5	5.1	6.4	0.7
Industrial	1.3	0.4	-0.3	0.7	1.2	0.1	0.0	0.8	0.9	-0.1	0.5	1.1
Miningo	-5.0	-4.0	-10.1	-6.1	-4.3	-4.1	-3.2	-4.8	-8.2	-6.5	-11.6	-8.5
Utilidad	2.5	3.3	-1.2	0.2	1.9	2.1	1.0	5.7	3.9	2.8	-0.5	-1.9
Comercio	3.9	2.2	-0.1	2.9	3.4	-0.7	1.3	3.0	0.0	3.0	1.5	-1.5
Manufacturas	2.8	1.1	3.6	2.7	2.6	2.0	0.7	1.5	1.1	1.9	5.2	2.0
Servicios	3.3	3.4	2.5	3.4	3.7	3.8	3.4	3.3	3.4	3.4	3.7	3.2
Comercio	4.9	3.0	3.0	4.4	5.1	4.3	3.7	3.3	1.3	2.6	3.4	2.4
Transportes	4.0	3.1	1.8	3.9	5.3	1.9	2.9	3.3	2.3	2.7	4.1	3.4
Info. en medios masivos	3.1	0.1	4.2	2.5	8.6	16.1	8.9	9.2	13.4	8.8	5.3	7.0
Financiero y de seguros	2.6	2.8	0.7	3.1	5.9	6.1	8.2	7.5	7.9	7.2	9.5	8.9
Inmobiliario y del alquiler	2.7	2.0	2.1	3.4	2.3	2.2	2.1	1.8	1.8	1.8	2.4	1.8
Resto	2.4	2.7	2.6	2.6	1.5	1.8	2.2	3.1	3.6	3.3	3.0	2.2

* En millones de pesos
Fuente: INEC



M	N	O	P	Q	R	S
ID_RECURSO	RECURSO	CLAVE_PROE	PROGRAMA	ID_RAMO	RAMO	INSTITUCION
2	Aportaciones I013		FONE Servici		33 Aportaciones	SECRETARÍA de
2	Aportaciones I012		FAFEF		33 Aportaciones	SECRETARÍA
3	Convenios K031		Proyectos de		9 Comunicacio	CARRETERAS
2	Aportaciones I002		FASSA		33 Aportaciones	SECRETARÍA
1	Subsidios U005		Seguro Popul		12 Salud	SERVICIOS D
3	Convenios K040		Proyectos de		9 Comunicacio	Secretaría de
1	Subsidios U006		Subsidios par		11 Educación P	UNIVERSIDA
1	Subsidios U006		Subsidios par		11 Educación P	Universidad
1	Subsidios U006		Subsidios par		11 Educación P	UNIVERSIDA
1	Subsidios U128		Proyectos de		23 Provisiones	Comisión Est.
1	Subsidios U128		Proyectos de		23 Provisiones	Comisión Est.
1	Subsidios U092		Fortalecimier		23 Provisiones	MUNICIPIO E
3	Convenios U006		Subsidios par		11 Educación P	UIED
1	Subsidios U130		Fortalecimier		23 Provisiones	Secretaría de
3	Convenios K031		Proyectos de		9 Comunicacio	CARRETERAS
3	Convenios K041		Sistema de Tr		9 Comunicacio	Secretaría de
1	Subsidios U005		Seguro Popul		12 Salud	SECRETARIA
1	Subsidios U005		Seguro Popul		12 Salud	SECRETARIA
1	Subsidios U022		Programas Ri		23 Provisiones	SECRETARIA
3	Convenios U006		Subsidios par		11 Educación P	UNIVERSIDA
2	Aportaciones I014		FONE Otros c		33 Aportaciones	Secretaría de
1	Subsidios U006		Subsidios par		11 Educación P	UNIVERSIDA
3	Convenios U012		Fortalecimier		12 Salud	SECOPE
1	Subsidios U057		Fondo Metro		23 Provisiones	Secretaría de
1	Subsidios U087		Fondo de Caj		23 Provisiones	Secretaría de
1	Subsidios U005		Seguro Popul		12 Salud	INSTITUTO D
1	Subsidios U006		Subsidios par		11 Educación P	Universidad
1	Subsidios U130		Fortalecimier		23 Provisiones	Tribunal Sup
3	Convenios K031		Proyectos de		9 Comunicacio	CARRETERAS

Example of structured vs unstructured data. On the left, data of the Gross Domestic Product that is public but difficult to process; on the right side, data on transfers of federal resources to sub-national governments in a structured format in rows and columns (XLSX and CSV) for agile processing.

3. Link between computer systems

When the information is already structured in storage systems, it is advisable to proceed to diagnose its quality and know its contents; design a proposal for the public release of the information (in case the format and structure are not usable) and perform tests to finally publish it.

In this situation, it is convenient to aspire to the highest levels of openness from the beginning of the process, in order to have an advantage in its publication. To learn more about the levels of data openness, [review the video and reference document](#) "What is open data?".

With these considerations on the identification of the source and format of the data, **the next step would be to start its debugging and cleaning.**

Check the next review document to learn about this process.



References

Kyocera. Difference between structured and unstructured data. Available at:
<https://smarterworkspaces.kyocera.es/blog/diferencia-datos-estructurados-no-estructurados/>

School of Data. Know the format of your data. Available at:
<https://es.schoolofdata.org/2013/10/31/conoce-el-formato-de-tus-datos/>

