

Marco de vinculación de datos abiertos aplicado al contexto de datos medioambientales

Juan Santiago Preisegger¹, Alejandro Greco, Ariel Pasini,

Marcos Boracchia, Patricia Pesado

Instituto de Investigación en Informática LIDI (III-LIDI)*

Facultad de Informática – Universidad Nacional de La Plata 50 y 120 La Plata Buenos Aires

*Centro Asociado Comisión de Investigaciones Científicas de la Pcia. de Bs. As. (CIC)

¹ Becario UNLP

{jspreisegger, apasini, marcosb, ppesado}@lidi.info.unlp.edu.ar
{alegreounlp}@gmail.com

Abstract. Los portales de datos abiertos ponen un conjunto muy importante de información a disposición de la comunidad. Los interesados en una temática en particular obtienen los datos sobre dicha temática desde diferentes portales, pero en el momento de querer procesarla en conjunto se dificulta debido a los diferentes criterios de cada portal para publicarlos. Para asistir en este proceso se desarrolló un marco de trabajo que facilita la vinculación de estos archivos a usuarios con poca experiencia técnica en el análisis de datos. Dentro del marco se utiliza una herramienta que permite aplicar diferentes operaciones sobre los archivos, generando gráficos en un panel de control, facilitando a los interesados el análisis de la información. El marco se aplicó a la temática del medioambiente, en particular, sobre la calidad del agua, del aire y la generación de energía.

Keywords: Datos Abiertos, Vinculación de Datos, Ingeniería de Software, Gobierno Abierto, Medioambiente.

1 Introducción

Una ciudad se nutre del comportamiento de sus ciudadanos. Los ciudadanos, mediante diferentes dispositivos, son capaces de registrar cada vez más información de las actividades que realizan. Hacer un uso inteligente de la información registrada por las autoridades gubernamentales para mejorar la vida de los ciudadanos es de gran aporte a la misma comunidad. Los aportes que se pueden lograr de los datos registrados no solo pueden provenir del área gubernamental, sino también de diferentes organismos o personas físicas, capaces de analizar la información, procesarla y proponer mejoras. Esto es un factor importante en este ciclo de mejora de una ciudad y es lo que la convierte en una “ciudad inteligente sostenible y participativa”. [1]

Para lograr una participación ciudadana en este tipo de procesos, los organismos ponen a disposición de su comunidad grandes volúmenes de datos abiertos, con el fin de que aquellos interesados en la temática puedan procesarlos y generar aportes en el proceso de mejora de la ciudad [2][3]. Pero al momento de acceder a los datos, aparecen diferencias técnicas como formatos de archivos, estructura del archivo, nombres de columnas, tipos de datos, magnitudes, etc., las cuales dificultan, y en algunos casos imposibilitan, el análisis de la información.

El marco de trabajo que se propone busca vincular datasets, obtenidos de portales abiertos, de una temática en particular y permitir un análisis en conjunto de estos datos, de una forma sencilla que no requiera de conocimientos técnicos avanzados. Dicho marco se basa en cinco pasos: 1) *Búsqueda*, 2) *Análisis preliminar*, 3) *Carga directa*, 4) *Normalización*, 5) *Vinculación*. Este proceso se apoya en la utilización de la herramienta Indimaker.

Indimaker tiene el potencial de procesar archivos de varios formatos, aplicar diferentes operaciones al contenido y lograr una vinculación entre los archivos, conformando un panel de control de indicadores que facilita la visualización de las operaciones realizadas sobre los archivos.

Para validar el proceso, el marco se aplicó a datos abiertos relacionados con el medioambiente, en particular a conjuntos de datos obtenidos de diferentes portales públicos sobre la calidad del aire, del agua y el uso de la energía.

En la segunda sección se presentan los conceptos de ciudades inteligentes sostenibles y datos abiertos. Luego, en la tercera sección, los conceptos generales sobre los tableros de control de indicadores y la herramienta Indimaker. En la cuarta sección se presenta el marco de trabajo de vinculación de datos abiertos. En la quinta, la aplicación del marco a datos relacionados con el medioambiente y, por último, en la sexta sección las conclusiones y trabajos futuros.

2 Las ciudades inteligentes sostenibles y los datos abiertos

Las ciudades se nutren de la participación de su comunidad. Los ciudadanos, constantemente, generan información que luego podrá ser utilizada en la toma de decisiones sobre el desarrollo de esa ciudad y, pasado un tiempo, influirá en la vida de esos ciudadanos. Lograr que los ciudadanos tengan acceso a los datos y se les permita participar en su análisis para construir el desarrollo de la ciudad, es un importante aporte a fin de convertir a una ciudad en una ciudad inteligente.

2.1 Ciudades inteligentes sostenibles

Generalmente, se relaciona el concepto de ciudades inteligentes con el uso de la tecnología en el desarrollo de las actividades de la ciudad, pero en realidad es mucho más que eso. Según [1] las ciudades inteligentes sostenibles representan las últimas etapas de progresión a través de ciudades digitales y ciudades inteligentes, consideradas como un proceso transformador continuo, basado en la colaboración y el compromiso de diferentes actores, construyendo diferentes capacidades (humanas,

técnicas e institucionales) de manera de mejorar la calidad de vida, proteger los recursos naturales y perseguir el desarrollo socioeconómico. La Unión Internacional de Telecomunicaciones (UIT) de la ONU, estableció una de las definiciones pioneras de ciudad inteligente sostenible: "Una ciudad inteligente sostenible es una ciudad innovadora que utiliza tecnologías de información y comunicación (TIC), y otros medios, para mejorar la calidad de vida, la eficiencia de la operación y los servicios urbanos y la competitividad, al tiempo que se garantiza que satisfaga las necesidades de las generaciones, presentes y futuras, con respecto a los aspectos económicos, sociales, ambientales y culturales".

2.2 Datos Abiertos

La sociedad presenta exigencias cada vez más elevadas hacia sus gobernantes y sus entes de gobierno. Entre estas exigencias, se incluye la transparencia y el manejo eficiente de los bienes públicos, la inclusión en la toma de decisiones y la colaboración con distintos sectores de la sociedad [4]. En función de estos requisitos, y con la ayuda de nuevas tecnologías, se generó una nueva forma de gobierno que incluye más al ciudadano, permitiéndole generar aportes a las políticas públicas y participar en la toma de decisiones. [5][6]

A través de la implantación del gobierno abierto se generó la apertura de los datos, la cual consiste en poner a disposición de la sociedad los datos de interés común de la ciudadanía para que, de cualquier forma, éstos puedan desarrollar una nueva idea o aplicación que entregue nuevos datos, conocimientos u otros servicios que el gobierno no es capaz de entregar. [3][7][8]

La variedad de los datos que se ponen a disposición de la sociedad es muy amplia y no solo se realiza desde las agencias gubernamentales, sino que, además, organismos internacionales, ONGs y otras organizaciones impulsan diversas medidas para poner a disposición de la sociedad cada vez más fuentes de datos, no solo referidas a las gestiones de gobierno, sino también relacionados a otros ámbitos, como el uso racional de los recursos y los cuidados del medioambiente.

3 Generación de tableros de control con datos abiertos

Una de las aplicaciones de los datos es la construcción de tableros de control, basados en indicadores generados con estos datos. Para la generación de estos tableros se deben analizar las fuentes de los datos, los datasets y poseer herramientas que nos permitan relacionar los datos, con el fin de construir información relevante que se transforme en un indicador que nos posibilite mejorar la toma de decisiones.

3.1 Fuentes de datos

Los nuevos paradigmas, gestados en los diferentes organismos, están coordinando acciones para mejorarle la calidad de vida a la sociedad, mediante la apertura de datos y las mejoras planteadas en la infraestructura de las ciudades. Se generaron diversas

aplicaciones y herramientas, desde múltiples sectores, para brindarles soporte y automatizar, o mejorar, el proceso de publicación, búsqueda y, en ocasiones, procesamiento de la información a las distintas partes de la sociedad. Entre estas herramientas se destacan los catálogos y portales de datos abiertos, mediante los cuales, los organismos publican los datos sobre diferentes aspectos de su ejercicio y del medioambiente en el cual operan. Por ejemplo, algunos países, provincias, municipios u organismos poseen portales donde unifican datos de las distintas regiones o temáticas en las que se especializan. Los más avanzados en el área, publican sus datos y los describen con esquemas de datos, para hacerlos más descriptivos en cuanto a su contenido.

Una herramienta que utiliza esta información para permitir llevar a cabo la búsqueda de los datos de manera más global es Google Dataset Search: un motor de búsqueda especializado en encontrar conjuntos de datos (datasets) almacenados en la web, a través de palabras clave, siempre y cuando estos utilicen etiquetas de conjuntos de datos schema.org o estructuras equivalentes representadas en el formato Data Catalog Vocabulary (DCAT).[9]

3.2 Datasets

El término dataset es un anglicismo. Su traducción al castellano sería “conjunto de datos” y es una colección de datos habitualmente tabulada, es decir, que corresponden a los contenidos de una única tabla de base de datos o una única matriz de datos, en donde cada columna de la tabla representa una variable en particular y cada fila representa un miembro determinado del conjunto de datos.

Los datasets son la pieza principal de los portales y catálogos de datos. Agrupan uno o más recursos de datos y, para su publicación, requieren una preparación previa para poder ser procesados y reutilizados por terceros. Según [10], esto incluye tres actividades: 1) *Documentación*: consiste en definir los metadatos que tendrá cada uno de los datasets a publicar. Los metadatos describen los aspectos básicos del dataset, y se utilizan para organizar, clasificar, relacionar y encontrar los datos necesarios (por ejemplo: título, descripción, institución, licencia, categoría, fecha de publicación, etc.). 2) *Estructuración*: consiste en la preparación del conjunto de datos a publicar en un formato estructurado, sin campos erróneos o vacíos, que permita su reutilización y procesamiento en cualquier software. 3) *Carga de datos*: consiste en su publicación en una plataforma, que permita la organización y fácil acceso por parte de quienes van a reutilizar los datos.

3.3 Tableros de control

Un indicador se puede definir como un dato, o conjunto de datos, que ayudan a medir objetivamente la evolución de un proceso o de una actividad, correspondiente a cualquier organización. Los indicadores se pueden organizar y relacionar de manera que conformen un tablero de control. Estos tableros permiten realizar un seguimiento y evaluación del proceso o actividad, de una forma más exhaustiva. Además,

generalmente, permiten ver la evolución de los mismos de forma gráfica facilitando la interpretación de los resultados.

3.4 Herramienta para la generación de tableros de control - Indimaker

IndiMaker es un sistema accesible desde cualquier navegador web, que permite la vinculación entre datasets, para la construcción de indicadores personalizados en tableros de control.

Los tableros de control se generan a partir de los distintos datasets cargados por el usuario en un primer momento, aunque posteriormente se permite agregar, quitar o realizar una combinación entre dos de estos, obteniendo un nuevo dataset listo para poder ser utilizado.

La herramienta permite importar datasets de diversos formatos (.xls, .xlsx, .xml, .ods y .csv), e incluso unir los que tengan columnas con contenido en común, para generar uno nuevo con información más enriquecida. A su vez, posee la capacidad de homogeneizar toda esta información y almacenarla en su base de datos, brindando la posibilidad al usuario de realizar operaciones entre los distintos datos almacenados para la construcción de indicadores.

Los indicadores deben ser generados de manera apropiada, sino, pueden generar información imprecisa, errónea o subjetiva constituyendo un obstáculo para el correcto análisis. Una de las grandes virtudes de esta herramienta es su sencillez a la hora de generarlos, ya que se encuentra diseñada para hacer uso de la información de una forma que resulte fácil para el usuario realizar operaciones complejas entre los diferentes datos, evitando a priori la generación de indicadores erróneos. Los indicadores generados servirán para elaborar una medida cuantitativa que tendrá significado para quien lo analice.

A su vez, la herramienta permite adaptarse a los diferentes dispositivos, ya sea una computadora de escritorio, una tablet o un smartphone, y además puede ser utilizada tanto en idioma español como en inglés.

Al acceder a la herramienta, se puede observar el menú para la gestión de los tableros de control pertenecientes al usuario, con sus fuentes de datos asignadas y una descripción que sirve para comprender el propósito de cada tablero. Para crear indicadores pertenecientes a un tablero, se debe indicar su nombre, tipo de indicador, usuario responsable, periodicidad de medición, descripción y niveles de referencia, además de la operación en la cual se basa. Las operaciones que pueden realizarse entre los datos incluyen sumas, restas, divisiones, porcentajes, agrupaciones de datos y varias comparaciones lógicas, incluida la posibilidad de utilizar expresiones regulares para usuarios más avanzados. La combinación de estas operaciones serán las que le permitan al usuario construir un indicador.

Los indicadores de cada tablero se pueden representar de varias formas, incluyendo diferentes gráficos, lo que permitirá alertar fácilmente al usuario sobre posibles desviaciones de los objetivos previamente establecidos.

4 Marco de vinculación de datos abiertos

Al analizar los datos de los diferentes catálogos o portales, aparecen problemas de incompatibilidad entre los formatos de los datos presentados por los diferentes proveedores. Dado este inconveniente, se comprendió que, para lograr una vinculación de datos exitosa, es de vital importancia analizar previamente las fuentes de datos y los formatos de cada una de ellas. En consecuencia, se avanzó en generar un marco de vinculación de datos de cinco pasos: 1) *Búsqueda*, 2) *Análisis preliminar*, 3) *Carga directa*, 4) *Normalización*, 5) *Vinculación*.

1. *Búsqueda*: Realizar la búsqueda de los datasets referidos al área de interés, mediante los portales o catálogos de datos de las organizaciones, o ayudándose con la herramienta Google Dataset Search.
2. *Análisis preliminar*: Analizar los datasets obtenidos, verificando que posean la información en un formato compatible con la herramienta y que todas las columnas posean encabezado. Además, comprobar el formato del contenido (completar filas, columnas, o trasponer filas por columnas en alguno de los datasets, etc.) para lograr una vinculación exitosa.
3. *Carga de archivos*: Cargar el archivo en Indimaker. Al cargar el archivo, la herramienta realizará una serie de verificaciones sobre el contenido. De ser validado, se procederá a la instancia de comparación. En el caso que detecte inconsistencias, se procederá a la instancia de normalización del contenido.
4. *Normalización*: El proceso de normalización puede realizarse de forma manual o automática (utilizando una herramienta externa), dependiendo del tamaño del archivo. En esta instancia se verifica:
 - a. Que los datos en las columnas de los datasets presenten el mismo tipo de dato.
 - b. Que los datos en las columnas posean algún valor.
 - c. Que no se exceda la cantidad de datos que puede poseer cada fila.
5. *Vinculación*: Unificar las tablas en base a un parámetro en común, en caso de ser necesario, para obtener más información y poder relacionar los datos para analizarlos de manera sencilla y directa.

5 Aplicación de la vinculación de datos abiertos de medioambiente

Los organismos gubernamentales y no gubernamentales ponen, a disposición de la comunidad, numerosos conjuntos de datos del medioambiente de la región geográfica a la que pertenecen. Al momento de analizar la información en su conjunto para obtener valores regionales, aparecen los problemas de incompatibilidad entre los formatos de los datos presentados por los diferentes organismos. Es de vital importancia analizar las fuentes de datos, y los formatos de cada una de ellas, para permitir realizar un proceso de vinculación de la información.

A continuación, se aplicarán los 5 pasos del marco de trabajo:

1. Búsqueda

Se seleccionó como caso de aplicación el medioambiente. Entre las distintas partes que componen la temática, se decidió acotarla a los aspectos que afectan a la sociedad día a día, como son la calidad del agua y del aire, bienes fundamentales para la vida diaria de la sociedad y el futuro de la misma. A su vez, se incluyó a la generación de energía, bien que se relaciona a menudo con la calidad de vida.

Las búsquedas se realizaron bajo las palabras: *Calidad de agua potable*, *Generación de energía*, *Calidad del aire*. En la tabla 1, a continuación, se puede observar la selección de los datasets obtenidos de los diferentes países.

<i>País</i>	<i>Agua</i>	<i>Aire</i>	<i>Energía</i>
<i>EEUU</i>	✓	✓	✓
<i>Chile</i>	✓	✓	✓
<i>Brasil</i>	✓	✓	✓
<i>Uruguay</i>	×	✓	✓
<i>Paraguay</i>	×	×	✓
<i>Bolivia</i>	×	×	×
<i>Perú</i>	×	×	×
<i>Colombia</i>	✓	✓	×
<i>Argentina</i>	✓	×	✓

Tabla 1. Datasets obtenidos por país

2. Análisis preliminar

A partir de los datos obtenidos se analizó su forma, los formatos disponibles y las estructuras de cada uno de ellos, de forma que, utilizando la herramienta propuesta, se los relacione para poder llevar a cabo un análisis mayor de la situación de las distintas regiones.

Datos de aire y agua

En el caso de la calidad del aire y la calidad del agua, fue posible visibilizar cierto estándar para analizar las magnitudes existentes en distintas características. Se pudo observar que se llevan a cabo casi las mismas pruebas sobre distintas muestras para analizar distintas características y determinar si se encuentran dentro de los márgenes saludables para el consumo humano.

Datos de la energía

En el caso de los datasets del área de energía, la diferencia entre los datos publicados por las distintas organizaciones es aún mayor. Se observó que, algunos países, publican simplemente el porcentaje anual generado y el tipo de obtención de energía en la cual se basan las distintas centrales de generación publicadas, otros países simplemente publican los porcentajes de sus fuentes de generación de energía, sin discernir entre las centrales que poseen, y otros publican las capacidades totales de las centrales, sus tipos y hasta las empresas propietarias de las mismas.

3. Carga directa

Con los datos recolectados se procedió a realizar la carga de manera directa, como punto de partida, para analizar si cumplía con el control que la herramienta realiza al momento de la carga de los datos. En este caso se comprobó que no todos los archivos poseían los datos de manera ordenada y sin errores, y solamente un 35% de los datasets pasaron este mínimo control, lo cual indica que más allá de la publicación de los datos realizada por los organismos, hace falta mayor normalización para su procesamiento. Si se analizan los datos por área, se puede observar que la progresión, en base a esa normalización, es: Agua: 80% de forma directa y un 20 % de forma indirecta; Energía: 66% Directa y 33% indirecta; Aire: 50% directa y 50% indirecta.

4. Normalización

En las tres áreas se pudo observar que, más allá de la estandarización en cuanto a ciertos datos publicados, los datos están en un formato muy “crudo” que imposibilita un procesamiento óptimo. En muchos casos poseen espacios vacíos, o incluso poseen más datos en las filas de los que se declara. En cada caso podemos ver:

Datos de aire y agua

Existen diferencias en la estructura de los datasets, las unidades utilizadas entre las distintas magnitudes que se analizan y cómo se almacenan estas en los datasets. Por ejemplo, la variación de columnas a filas o la división de un campo en varios. Esto dificulta la interrelación entre los distintos datasets de estas temáticas.

Datos de la energía

En este caso, se observó que varía mucho la unidad de medición utilizada en los distintos datasets, lo cual imposibilita una comparación lineal de los datos. A su vez, por la gran diferencia en cuanto a los datos publicados por cada organización, es necesario llevar a cabo una selección de ciertos campos en común entre los distintos datasets, sumándolos en base a su tipo, para poder visibilizarlos de manera tipificada. Esto se dificulta debido a la dispersión existente.

5. Vinculación

Una vez normalizados los datos, se puede realizar la vinculación de distintos datasets que posean datos en común, permitiendo un análisis regional de la información. Con estos datasets, la herramienta nos permite seleccionar columnas para realizar operaciones y obtener valores con los que se pueden generar diversos gráficos que permiten un análisis lineal de magnitudes de interés. A continuación, se puede visibilizar el análisis que se realizó en cada uno de los conjuntos de datos analizados.

Datos del agua

Datos de la energía

Entre la información recolectada sobre la producción de energía se seleccionaron, a modo de ejemplo, los datos de Paraguay y la ciudad de Nueva York. Más allá de sus diferencias en cuanto a dimensión y a población, se realizó la comparación con estos datos para poder observar las diferencias en cuanto a los valores de la producción energética, en GWh. La información publicada en este caso era muy similar, aunque Nueva York especificaba, además, las fuentes de producción de energía. Con esta información, se generaron los gráficos de líneas que se pueden observar a continuación, en la Imagen 3, donde se puede apreciar que la ciudad de Nueva York produjo, en el período analizado, más energía que todo Paraguay.

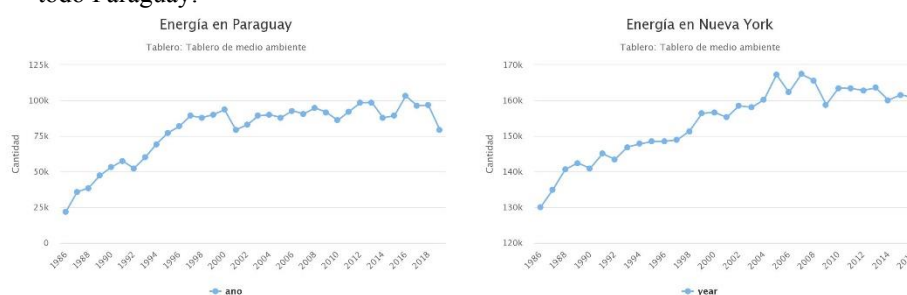


Imagen 3. Producción de energía por año en Paraguay y en Nueva York

6 Conclusión y trabajos futuros

A lo largo del artículo se introdujeron los conceptos básicos de Ciudades Inteligentes Sostenibles, Datos Abiertos, Fuentes de Datos, Datasets y Tableros de Control. Se hizo hincapié en cómo pueden relacionarse estos conceptos para generar un mayor nivel de información para la sociedad, a partir de la vinculación de los datos que la misma dispone y en base al uso de la herramienta Indimaker.

En función de esta herramienta se generó un marco de trabajo, compuesto por cinco etapas, para la vinculación de los datos puestos a disposición, por diversas organizaciones, de manera de poder compararlos sencillamente.

Como caso de estudio, para la aplicación del marco, se descargaron fuentes de datos abiertos de 10 países sobre la calidad del agua, calidad del aire y la producción energética, obteniendo información en 8 de estos portales de datos.

La información obtenida se procesó, según el marco de trabajo, normalizando la información en los casos que fue necesario y, luego, utilizando la herramienta Indimaker, la cual brindó la posibilidad de realizar operaciones sobre los datos para generar indicadores de forma sencilla. Fue notable la diferencia que se encontró en la información publicada sobre la misma temática, pero, con la aplicación del marco propuesto, la herramienta pudo procesar la información generando indicadores de

manera sencilla para el usuario y diversos gráficos que facilitaron el análisis de cada área. Cabe destacar que el modelo puede extenderse a cualquier área de interés.

A futuro, se espera vincular la herramienta con APIs de diferentes portales de datos brindadas por diversas organizaciones y ampliar el abanico de operaciones que la herramienta brinda, incluyendo la posibilidad de visualizaciones de diversos indicadores en un mismo gráfico.

7 Agradecimientos

Project co-funded by the Erasmus+ Programme of the European Union. Grant no: 598273-EPP-1-2018-1-AT-EPPKA2-CBHE-JP.

8 Bibliografía

- [1] E. Estevez, N. V. Lopes, and T. Janowski, “Smart Sustainable Cities - Reconnaissance Study. Operating Unit ON Policy-Driven. Electronic Governance,” *United Nations University, Canada*, 2017. .
- [2] S. A. Chun, S. Shulman, R. Sandoval, and E. Hovy, “Government 2.0: Making connections between citizens, data and government,” *Inf. Polity*, vol. 15, no. 1–2, pp. 1–9, 2010.
- [3] G. Concha and A. Naser, “Datos abiertos: Un nuevo desafío para los gobiernos de la región,” *Inst. Latinoam. y del Caribe Planif. Económica y Soc.*, 2012.
- [4] C. Calderón and S. Lorenzo, *Open Government. Gobierno Abierto*. 2010.
- [5] A. Naser, Á. Ramírez-Alujas, and D. R. Editores, *Desde el gobierno abierto al Estado abierto en America Latina y el Caribe: Planificación para el Desarrollo*. 2017.
- [6] J. R. Gil-García and J. I. Criado, *Las Tecnologías de Información y Comunicación en las Administraciones Públicas Contemporáneas*. 2017.
- [7] A. Pasini, J. S. Preisegger, and P. Pesado, “Modelos de evaluación de gobiernos abiertos , aplicado a los municipios de la provincia de Buenos Aires,” *XXIV Congr. Argentino Ciencias la Comput.*, vol. XXIV, pp. 0–10, 2018.
- [8] A. Pasini, J. S. Preisegger, and P. Pesado, “Open Government Assessment Models Applied to Province’s Capital Cities in Argentina and Municipalities in the Province of Buenos Aires,” in *Communications in Computer and Information Science*, 2019, vol. 995, pp. 355–366.
- [9] N. Noy, M. Burgess, and D. Brickley, “Google dataset search: Building a search engine for datasets in an open web ecosystem,” in *The Web Conference 2019 - Proceedings of the World Wide Web Conference, WWW 2019*, 2019, pp. 1365–1375.
- [10] A. Naser and A. Ramirez, “Plan de gobierno abierto. Una hoja de ruta para los Gobiernos de la Región,” *CEPAL - Manuales*, vol. 81, p. 80, 2017.