

Visualización Avanzada y Automatización del Análisis de
Datos

Tema 3. Preprocesamiento y estructuración de datos para visualización

Índice

Esquema

Ideas clave

- 3.1. Introducción y objetivos
- 3.2. Principios de organización de datos
- 3.3. Limpieza y transformación en herramientas visuales
- 3.4. Procesos de preprocesamiento con lenguajes de programación
- 3.5. Integración de múltiples fuentes
- 3.6. Preparación de datos para dashboards
- 3.7. Resumen y conclusiones
- 3.8. Referencias bibliográficas

A fondo

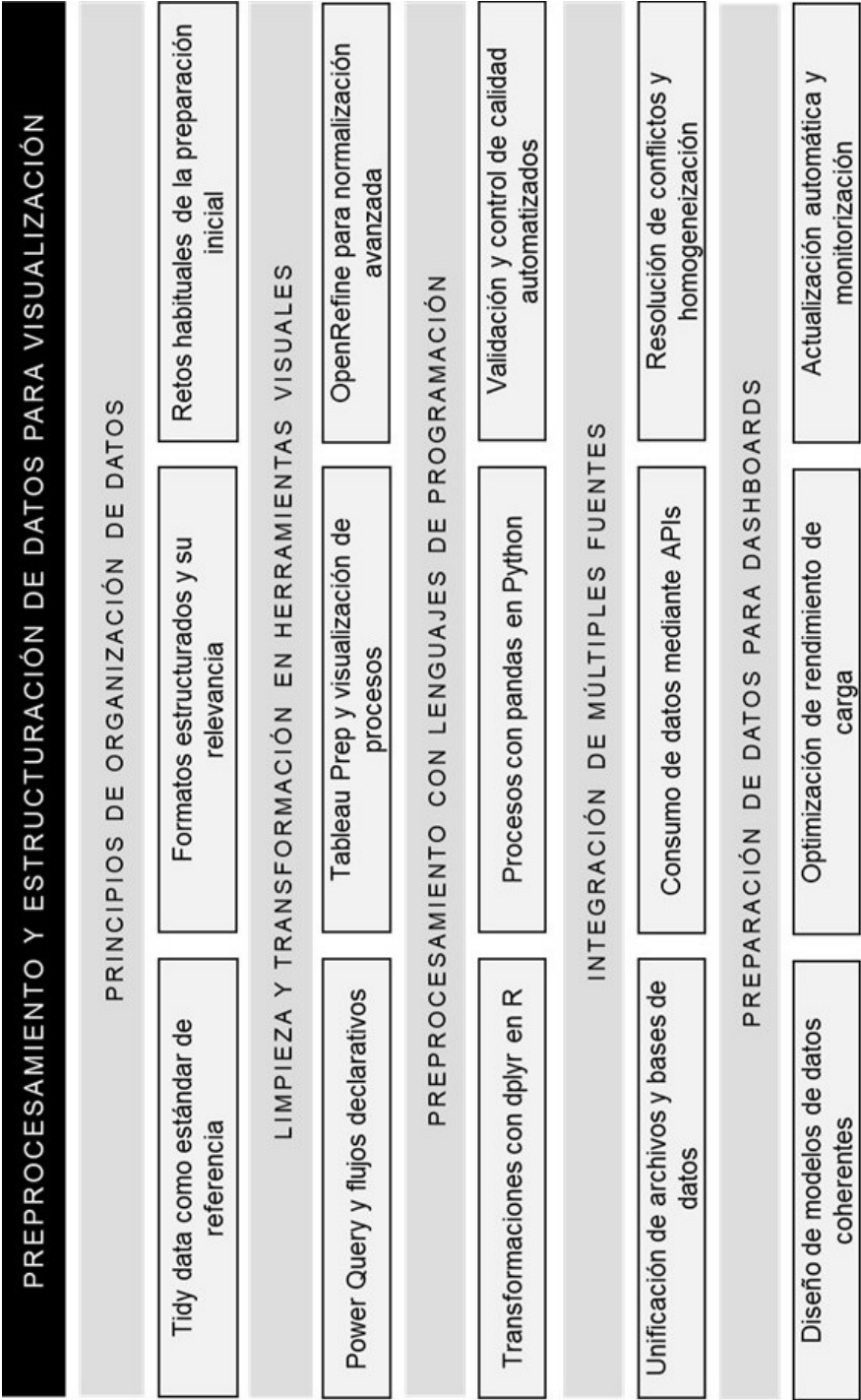
The Art of Data Cleaning: Transforming Messy Data into Structured Insights

Guide To Data Cleaning: Definition, Benefits, Components, And How To Clean Your Data

OpenRefine Tutorial – University of Washington

Tutorial de OpenRefine: limpieza y transformación de datos abiertos

Test



3.1. Introducción y objetivos

Pocos profesionales se detienen a pensar que más del 70 % del tiempo invertido en proyectos de análisis visual no se dedica a diseñar gráficos, sino a preparar los datos. Este porcentaje, que sorprende a quienes se inician en la disciplina, refleja una verdad ineludible: la calidad de cualquier visualización depende directamente de la solidez del preprocesamiento. Una base de datos ordenada, limpia y bien estructurada es la diferencia entre un *dashboard* que aporta valor y otro que conduce a conclusiones erróneas.

A diario, las organizaciones se enfrentan al reto de integrar información proveniente de hojas de cálculo dispersas, sistemas transaccionales heredados y servicios externos que generan flujos de datos en tiempo real. En este contexto, la limpieza y transformación de datos no solo es una cuestión técnica: es un proceso estratégico que exige criterio, precisión y un conocimiento profundo de las herramientas. Power BI, Tableau Prep, R y Python se han consolidado como entornos clave que permiten afrontar esta complejidad con metodologías reproducibles y auditables.

Este tema proporciona una guía detallada sobre los principios y técnicas que sustentan el preprocesamiento y la estructuración de datos orientados a la visualización avanzada. Se revisarán los fundamentos del enfoque *tidy data*, los procesos de transformación en plataformas visuales y lenguajes de programación, la integración de fuentes heterogéneas y las estrategias de optimización que facilitan la creación de modelos de datos robustos. Gracias a estos conocimientos, el alumnado podrá asumir con solvencia los desafíos que plantea la preparación de datos en escenarios reales.

Al finalizar este tema, el alumnado será capaz de:

- ▶ Comprender los principios de organización y estructuración de datos (*tidy data*) y su impacto en la calidad del análisis visual.
- ▶ Aplicar técnicas de limpieza y transformación de datos mediante Power BI, Tableau Prep, R y Python, evaluando las ventajas y limitaciones de cada entorno.
- ▶ Integrar y unificar datos procedentes de diversas fuentes, resolviendo inconsistencias y optimizando su preparación para su uso en *dashboards* interactivos y automatizados.

3.2. Principios de organización de datos

La organización de los datos es el punto de partida imprescindible para cualquier proyecto de visualización. Antes de pensar en gráficos, colores o interacciones, es necesario garantizar que la información sigue un esquema lógico, estandarizado y libre de ambigüedades. Este proceso no solo facilita el análisis, sino que permite mantener la coherencia y la reproducibilidad del trabajo a lo largo del tiempo.

A lo largo de este apartado se explorarán los fundamentos que sustentan el concepto de *tidy data*, los estándares que definen un conjunto de datos bien estructurado y los retos más frecuentes a los que se enfrentan los profesionales cuando preparan información para su explotación visual. Comprender estos principios es esencial para que las herramientas de visualización puedan operar con fiabilidad y eficiencia.

El concepto de tidy data y su importancia

El término *tidy data* hace referencia a un conjunto de principios que definen cómo deben organizarse los datos para que sean comprensibles y procesables por herramientas de análisis y visualización. Esta idea fue popularizada por Hadley Wickham, y establece que cada variable debe formar una columna, cada observación debe ocupar una fila y cada tipo de unidad observacional debe almacenarse en una tabla independiente. Aplicar esta estructura simplifica las transformaciones, reduce los errores y facilita la integración de distintas fuentes.

En la práctica, trabajar con datos ordenados implica identificar inconsistencias y reestructurar tablas que a menudo contienen mezclas de variables y etiquetas. Por ejemplo, es habitual encontrar hojas de cálculo donde los nombres de los meses aparecen como cabeceras de columna en lugar de valores en una variable «mes» o donde los totales se mezclan con los registros detallados. Estas situaciones dificultan la automatización de procesos y multiplican el riesgo de interpretación errónea.

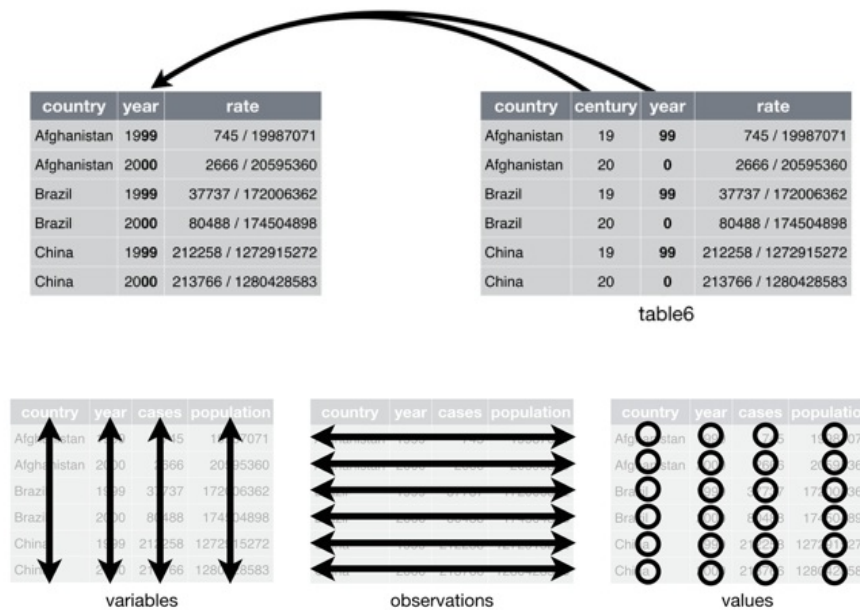


Figura 1. Ejemplos del concepto de datos ordenados (tidy data). Las imágenes muestran cómo reorganizar tablas para que cada variable ocupe su propia columna, cada observación una fila y cada valor una celda única, facilitando así la transformación, el análisis y la visualización consistentes. Fuente: <https://r4ds.had.co.nz/tidy-data.html>.

La adopción del enfoque *tidy data* permite establecer un lenguaje común entre herramientas, desarrolladores y analistas. Cuando los datos se organizan siguiendo estos principios, la creación de gráficos, modelos predictivos o cuadros de mando se convierte en un proceso mucho más ágil y menos propenso a errores. Por esta razón, entender y aplicar este concepto constituye uno de los pilares de la visualización avanzada.

Estándares y formatos de datos estructurados

Más allá del concepto de *tidy data*, existen estándares ampliamente aceptados que definen cómo debe representarse la información en distintos contextos. Entre ellos destacan los formatos tabulares en CSV, los esquemas de bases de datos relacionales y las estructuras jerárquicas como JSON o XML. Cada uno de estos formatos presenta ventajas y limitaciones que deben evaluarse en función del proyecto y la herramienta que se vaya a utilizar.

Los formatos tabulares, por ejemplo, son ideales para datos transaccionales o de series temporales que se consumen en aplicaciones de BI. Sin embargo, en entornos donde la información proviene de APIs o flujos en tiempo real, formatos como JSON permiten manejar estructuras más complejas y anidadas. La clave está en seleccionar el esquema que minimice las transformaciones necesarias y garantice la integridad de los datos durante su ciclo de vida.

El conocimiento de estos estándares resulta especialmente relevante cuando se integran fuentes heterogéneas. La unificación de datos exige establecer correspondencias entre campos, definir claves de relación y resolver conflictos semánticos. Por ello, el profesional de visualización no solo debe manejar conceptos estadísticos, sino también entender los fundamentos de modelado de datos y las implicaciones de cada formato en la fase de preprocesamiento.

Desafíos habituales en la preparación de datos

La preparación de datos conlleva retos que van mucho más allá de la limpieza superficial de registros. Uno de los problemas más frecuentes es la presencia de valores faltantes, que pueden deberse a errores en la recogida de información o a la inexistencia de datos en ciertos casos. La forma en que se gestionan estos vacíos tiene un impacto directo sobre la interpretación: imputar valores, eliminarlos o dejarlos explícitos son decisiones que deben tomarse con criterio analítico.

Otro desafío habitual es la presencia de duplicados y registros inconsistentes. Cuando un mismo evento se registra en más de un sistema, es común que aparezcan discrepancias en fechas, cantidades o etiquetas. La «desduplicación» (eliminación de los duplicados) y la normalización son procesos clave para asegurar que cada observación sea única y que las variables mantengan una codificación coherente a lo largo de todo el conjunto de datos.

Por último, los datos pueden presentar errores semánticos que no son evidentes a simple vista. Por ejemplo, una variable numérica que mezcla diferentes unidades de medida, o un campo categórico donde la misma categoría se nombra de formas distintas. Estos problemas exigen un conocimiento profundo del dominio de negocio y un enfoque sistemático de validación. Identificar y resolver estos retos de forma temprana reduce significativamente el esfuerzo necesario en fases posteriores de análisis y visualización.

3.3. Limpieza y transformación en herramientas visuales

La proliferación de plataformas de análisis visual ha traído consigo un conjunto de herramientas orientadas a simplificar la limpieza y transformación de datos mediante entornos gráficos. Estas soluciones permiten a profesionales de distintos perfiles preparar información sin necesidad de recurrir exclusivamente a código, reduciendo las barreras de entrada y aumentando la trazabilidad de los procesos. Al mismo tiempo, ofrecen potentes opciones de automatización y documentación que aseguran la consistencia en proyectos complejos.

Este apartado presenta una visión comparada de los entornos más utilizados para el preprocesamiento visual: Power BI, Tableau Prep, OpenRefine y la metodología de evaluación de flujos de transformación. Conocer sus fortalezas y limitaciones, facilita seleccionar la herramienta adecuada según el volumen de datos, el tipo de usuario y el grado de automatización requerido.

Procesos en Power BI (Power Query)

Power Query es el entorno de transformación de datos integrado en Power BI. Su principal fortaleza radica en la combinación de una interfaz intuitiva con un motor de procesamiento altamente optimizado, capaz de gestionar grandes volúmenes de datos de múltiples orígenes. Cada paso de limpieza, —como filtrado, agrupación, conversión de tipos o eliminación de duplicados— se almacena como una transformación declarativa que puede revisarse y modificarse en cualquier momento.

Una de las ventajas diferenciales de Power Query es su capacidad de conexión con una amplia variedad de fuentes: desde archivos locales y bases de datos, hasta servicios en la nube y APIs. Cada consulta se define mediante un lenguaje específico, M, que permite personalizar procesos complejos más allá de las opciones de la interfaz. Esta flexibilidad convierte Power Query en una solución adecuada tanto para perfiles técnicos como para usuarios de negocio que requieran autonomía.

La trazabilidad es otro aspecto clave de Power Query. Todas las transformaciones quedan documentadas en una secuencia visual que facilita la auditoría y el mantenimiento. Además, la integración con el modelo de datos de Power BI permite reutilizar las tablas transformadas de forma directa en informes y dashboards interactivos, acelerando el ciclo de desarrollo de proyectos de BI.

Procesos en Tableau Prep

Tableau Prep es la propuesta de Tableau para la limpieza y transformación de datos. Su enfoque está orientado a la claridad visual: cada paso se representa mediante nodos en un flujo interactivo que muestra de forma inmediata cómo se modifica la información. Este modelo permite explorar datos, identificar errores y transformar registros con un nivel de transparencia especialmente valorado en entornos colaborativos.

El funcionamiento de Tableau Prep se basa en tres operaciones principales: conexión a datos, transformación y salida. Entre las transformaciones más comunes se encuentran la unión de tablas, la agregación de registros, el pivotado y la creación de campos calculados. Cada acción se añade de forma secuencial y los resultados se previsualizan en tiempo real, facilitando la detección temprana de problemas de calidad.

Una de las características más apreciadas de Tableau Prep es su integración nativa con Tableau Desktop. Los flujos de datos preparados pueden publicarse como extractos reutilizables en múltiples proyectos de visualización. Además, la

funcionalidad de programación mediante Tableau Prep Conductor permite automatizar la ejecución periódica de procesos de limpieza, garantizando que los datos estén siempre actualizados.

Procesos con OpenRefine

OpenRefine es una herramienta de código abierto orientada a la limpieza masiva y estandarización de datos heterogéneos. Aunque su interfaz recuerda a una hoja de cálculo, su filosofía de uso es muy distinta: cada operación se aplica de forma estructurada sobre columnas completas, permitiendo procesar grandes volúmenes de información de manera eficiente. OpenRefine destaca especialmente en tareas de exploración, detección de inconsistencias y normalización de valores.

Entre las funciones más potentes de OpenRefine se encuentran el agrupamiento por similitud (clustering) y la transformación mediante expresiones regulares. Estas capacidades permiten identificar variantes de un mismo valor —por ejemplo, nombres de clientes escritos de diferentes maneras— y unificarlos de manera semiautomática. Asimismo, la herramienta facilita importar datos desde APIs o exportar resultados a múltiples formatos, lo que incrementa su versatilidad.

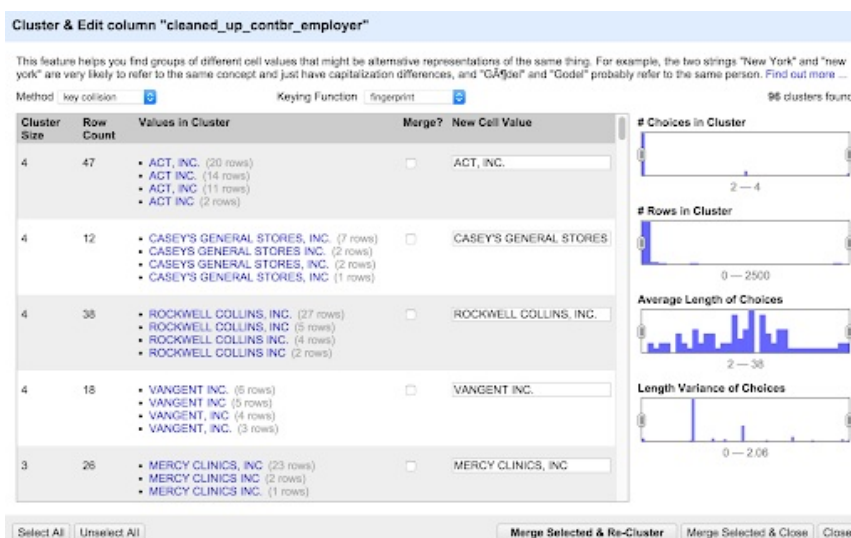


Figura 2. Vista del proceso de agrupamiento por similitud en OpenRefine, que permite identificar y unificar variantes de un mismo valor en una columna de datos. Fuente: <http://www.padjo.org/tutorials/open-refine/clustering/>.

La transparencia de OpenRefine es otro punto fuerte. Cada paso de limpieza queda registrado en un historial que puede exportarse como script JSON, lo que facilita reproducir el flujo de trabajo en otros proyectos o compartirlo con colaboradores. Por su carácter gratuito y su comunidad activa, OpenRefine se ha consolidado como un recurso accesible y eficaz en la etapa de preparación de datos, complementando otras herramientas más orientadas a la visualización.

Comparativa de flujos y mejores prácticas

Elegir entre Power Query, Tableau Prep y OpenRefine depende en gran medida del contexto del proyecto y de las competencias técnicas del equipo. Mientras Power Query sobresale en entornos empresariales con integración nativa en Power BI, Tableau Prep destaca por su transparencia visual y facilidad de uso, y OpenRefine aporta un nivel de control detallado sobre la limpieza y normalización que resulta difícil de igualar en otras soluciones gráficas.

Una buena práctica consiste en documentar claramente el flujo de transformación, describiendo cada paso con nombres significativos y anotaciones que expliquen el propósito de cada operación. Este enfoque no solo facilita la auditoría, sino que acelera el mantenimiento cuando los datos cambian o se incorporan nuevas fuentes. Además, conviene definir criterios de validación de calidad en cada etapa: por ejemplo, comprobar la ausencia de valores nulos o la coherencia de las claves primarias antes de continuar con la integración.

La combinación de herramientas también es una estrategia frecuente en proyectos complejos. Es habitual utilizar OpenRefine para la limpieza inicial, y posteriormente Power Query o Tableau Prep para la transformación y consolidación final. Esta sinergia permite aprovechar las fortalezas específicas de cada plataforma y construir procesos de preprocesamiento sólidos, replicables y fáciles de mantener.

3.4. Procesos de preprocesamiento con lenguajes de programación

La limpieza y transformación de datos mediante lenguajes de programación aporta un nivel de flexibilidad y control difícil de igualar con herramientas exclusivamente visuales. R y Python son hoy en día los entornos más empleados en proyectos de análisis avanzado y visualización, gracias a su capacidad de manejar grandes volúmenes de datos, integrar múltiples fuentes y automatizar procesos complejos con un alto grado de reproducibilidad. La programación permite crear flujos personalizados que se adaptan con precisión a los requisitos de cada proyecto.

El aprendizaje de estas técnicas no se limita al ámbito técnico: también proporciona al profesional de BI una mayor autonomía y capacidad de innovación. Dominar librerías como *dplyr* en R o *pandas* en Python no solo facilita la preparación de datos, sino que sienta las bases para tareas posteriores de modelización y generación de visualizaciones interactivas. Este apartado revisa las principales metodologías y buenas prácticas para trabajar con ambos lenguajes en el preprocesamiento.

Limpieza y transformación con R (tidyverse, dplyr)

R se ha consolidado como un lenguaje de referencia en la limpieza y estructuración de datos, especialmente a través del ecosistema tidyverse. Este conjunto de paquetes —entre los que destacan *dplyr*, *tidyr* y *readr*— proporciona una sintaxis coherente y expresiva que facilita tareas habituales como la filtración de registros, la creación de variables derivadas, la agregación y la normalización de datos. Una de las principales ventajas es la claridad de los flujos de trabajo, que pueden leerse casi como un lenguaje natural.

Dplyr es la piedra angular de estos procesos. Mediante funciones como `filter()`, `mutate()`, `summarise()` y `group by()`, es posible transformar conjuntos de datos de manera secuencial y reproducible. Por ejemplo, un flujo típico puede consistir en

importar un archivo CSV, eliminar registros incompletos, crear variables calculadas y agrupar resultados por categorías. Cada paso se encadena mediante el operador de tubería `>`, que aporta legibilidad y facilita el mantenimiento.

Otra característica destacada de R es su capacidad de integración con otros entornos. Es posible consumir datos de bases relacionales mediante *DBI* o *odbc*, trabajar con APIs usando *httr* y generar informes reproducibles en RMarkdown. Esta combinación convierte a R en una solución potente para proyectos en los que la limpieza de datos se vincula directamente con la generación de visualizaciones y la documentación automatizada.

Limpieza y transformación con Python (pandas)

Python es hoy el lenguaje más popular en el ecosistema de análisis de datos, gracias a su versatilidad y a la madurez de sus librerías. *pandas* es la herramienta más empleada para la manipulación de datos tabulares y ofrece una amplia gama de métodos que permiten importar, limpiar y transformar conjuntos de datos con rapidez. Su estructura basada en *DataFrames* facilita la transición conceptual desde entornos como Excel o SQL hacia flujos de trabajo programados.

Las operaciones más habituales con *pandas* incluyen la eliminación de duplicados, la imputación de valores faltantes, el cambio de tipos de datos, la normalización de variables y la creación de nuevas columnas mediante expresiones personalizadas. La sintaxis, aunque concisa, requiere comprender bien conceptos clave como los índices, la alineación de datos y las funciones de agregación. Una práctica recomendada es documentar cada transformación mediante comentarios y mantener scripts modularizados.

Python también destaca por su capacidad de integrarse con un ecosistema más amplio de librerías. Por ejemplo, *numpy* permite realizar operaciones numéricas de alto rendimiento, *sqlalchemy* facilita la conexión con bases de datos relacionales y

requests o *aihttp* permiten consumir datos de APIs externas. Este grado de flexibilidad convierte a Python en una elección excelente para proyectos que requieren tanto preprocesamiento como análisis estadístico y desarrollo de visualizaciones interactivas con librerías como *matplotlib* o *plotly*.

Validación y control de calidad de los datos

El control de calidad es un componente fundamental en cualquier flujo de preprocesamiento. Tanto en R como en Python, existen metodologías y funciones específicas para validar la integridad de los datos antes de que pasen a fases posteriores de análisis o visualización. La validación permite detectar incoherencias, asegurar que se cumplen las reglas de negocio y documentar los criterios aplicados en cada paso.

Una práctica habitual es definir reglas explícitas que comprueben la presencia de valores nulos en campos clave, la unicidad de las claves primarias o la coherencia de los rangos numéricos. En R, paquetes como *assertr* o *validate* facilitan la creación de estas reglas y su aplicación sistemática. En Python, librerías como *pandera* permiten definir esquemas de validación sobre los *DataFrames* de *pandas*, generando informes detallados sobre cada control.

La documentación de estos controles y su automatización forman parte de las mejores prácticas profesionales. Registrar los criterios aplicados y los resultados de cada validación permite demostrar la calidad del proceso y agiliza la resolución de incidencias. Además, incorporar la validación en el flujo de transformación incrementa la fiabilidad de las visualizaciones y refuerza la confianza de los usuarios finales en la información presentada.

3.5. Integración de múltiples fuentes

La integración de datos procedentes de fuentes heterogéneas es una de las etapas más exigentes del preprocesamiento. Las organizaciones trabajan con sistemas de información diversos: aplicaciones de gestión empresarial, plataformas en la nube, archivos locales, y servicios de datos en tiempo real. Cada origen utiliza formatos, convenciones y estructuras distintas que deben unificarse antes de generar un modelo de datos coherente y útil para la visualización.

Este proceso no se limita a importar registros: implica definir relaciones, resolver conflictos de nomenclatura, homogeneizar tipologías, y garantizar la integridad referencial. La integración eficaz requiere tanto conocimiento técnico de las herramientas como criterio analítico para determinar qué información es relevante y cómo debe vincularse. A continuación, se describen los principales aspectos que condicionan esta fase crítica del preprocesamiento.

Integración de archivos locales y bases de datos

El punto de partida más habitual en proyectos de análisis es la combinación de archivos planos (CSV, Excel) con datos almacenados en bases relacionales. Mientras los primeros ofrecen facilidad de acceso y flexibilidad en la preparación inicial, las bases de datos proporcionan mayor robustez, escalabilidad y control sobre la consistencia de la información. Integrar ambos orígenes requiere planificar cuidadosamente los procesos de extracción y conexión.

Herramientas como Power Query en Power BI y Tableau Prep ofrecen conectores nativos que permiten importar datos desde sistemas como SQL Server, MySQL u Oracle, así como desde carpetas locales. Una buena práctica es definir criterios de actualización automática que mantengan los datos sincronizados con las fuentes originales, evitando duplicidades o desfases temporales. Además, conviene documentar la procedencia de cada tabla o vista, ya que esta trazabilidad facilita la auditoría y la resolución de incidencias.

Un aspecto clave es la definición de claves de relación entre tablas. Identificar los campos que permitan vincular registros —por ejemplo, códigos de producto, identificadores de cliente o fechas— es imprescindible para consolidar la información y construir un modelo de datos que pueda explotarse con confianza. La consistencia de estas relaciones se verifica mediante reglas de validación que comprueban la ausencia de registros huérfanos y la unicidad de los identificadores.

Consumo de datos desde APIs

Cada vez más proyectos incorporan datos que provienen de servicios web o plataformas en la nube que exponen APIs RESTful. Este enfoque permite acceder a información en tiempo real, como métricas de redes sociales, tipos de cambio, catálogos de productos o series temporales de sensores. La integración de estos datos requiere comprender los formatos de intercambio —habitualmente JSON o XML— y desarrollar procesos de extracción y normalización adaptados a sus particularidades.

En entornos visuales como Power BI, los conectores web permiten consumir datos de APIs mediante peticiones GET o POST, configurando parámetros dinámicos y autenticación, si es necesario. En R y Python, paquetes como *httr*, *requests* o *aiohttp* ofrecen funciones para realizar estas conexiones y transformar la respuesta en estructuras tabulares que puedan combinarse con otras fuentes. La programación resulta especialmente útil cuando es preciso automatizar la paginación o manejar límites de volumen en las respuestas.

Una recomendación habitual es almacenar temporalmente las descargas de datos en repositorios intermedios, como bases SQL o sistemas de almacenamiento en la nube, para controlar versiones y facilitar la trazabilidad. De este modo, si una API cambia su esquema o limita su disponibilidad, se dispone de un histórico que asegura la continuidad del análisis. Además, es importante establecer procesos de control de calidad que validen la integridad y coherencia de los datos obtenidos.

Resolución de conflictos y homogeneización

Una de las mayores dificultades al integrar fuentes heterogéneas es la aparición de conflictos entre campos con significados equivalentes pero formatos diferentes. Por ejemplo, una base de datos puede almacenar fechas en formato ISO (AAAA-MM-DD), mientras que un archivo Excel utiliza notaciones regionales con día y mes invertidos. Del mismo modo, las categorías pueden variar en nomenclatura o codificación, lo que obliga a definir reglas de normalización explícitas.

La homogeneización de datos implica transformar las variables para que utilicen un mismo estándar. Este proceso puede incluir conversiones de tipos, mapeo de categorías a un diccionario unificado, estandarización de unidades de medida o eliminación de caracteres especiales en identificadores. Herramientas como Power Query y Tableau Prep permiten definir estas transformaciones mediante pasos secuenciales, mientras que en R o Python se desarrollan scripts que aplican las reglas de forma programada.

Otro aspecto clave es la resolución de duplicidades y solapamientos. Cuando la misma información procede de dos sistemas distintos, es frecuente encontrar registros redundantes con pequeñas discrepancias. La definición de criterios de prioridad —por ejemplo, tomar como referencia la fuente más actualizada o la que tiene mayor confiabilidad— es una práctica que debe consensuarse con los responsables del dominio de negocio. Documentar cada decisión garantiza la transparencia y facilita la gestión de futuras integraciones.

3.6. Preparación de datos para dashboards

Una vez que los datos han sido limpiados, transformados e integrados, es necesario adaptarlos a los requisitos específicos de los dashboards que se van a construir. Esta fase implica diseñar estructuras que faciliten la exploración visual, optimicen el rendimiento de carga, y permitan actualizar la información de manera eficiente. La calidad de un cuadro de mando no depende únicamente de su diseño gráfico: su capacidad de respuesta y fiabilidad, comienzan en esta etapa de modelado.

El diseño del modelo de datos adecuado requiere comprender las necesidades del usuario final, definir relaciones entre tablas que soporten los indicadores clave, y planificar estrategias de actualización que aseguren que el contenido permanece vigente. Este apartado presenta los principales elementos que intervienen en la preparación de datos orientada a la visualización avanzada.

Diseño de tablas modelo y relaciones

El punto de partida en la preparación de datos es la definición de las tablas modelo que soportarán las visualizaciones. Estas tablas deben estar estructuradas de manera que reflejen con claridad las entidades del negocio —clientes, productos, operaciones— y permitan responder a las preguntas más habituales que plantearán los usuarios. La claridad en la definición de las relaciones entre tablas es fundamental para garantizar la consistencia y evitar duplicidades.

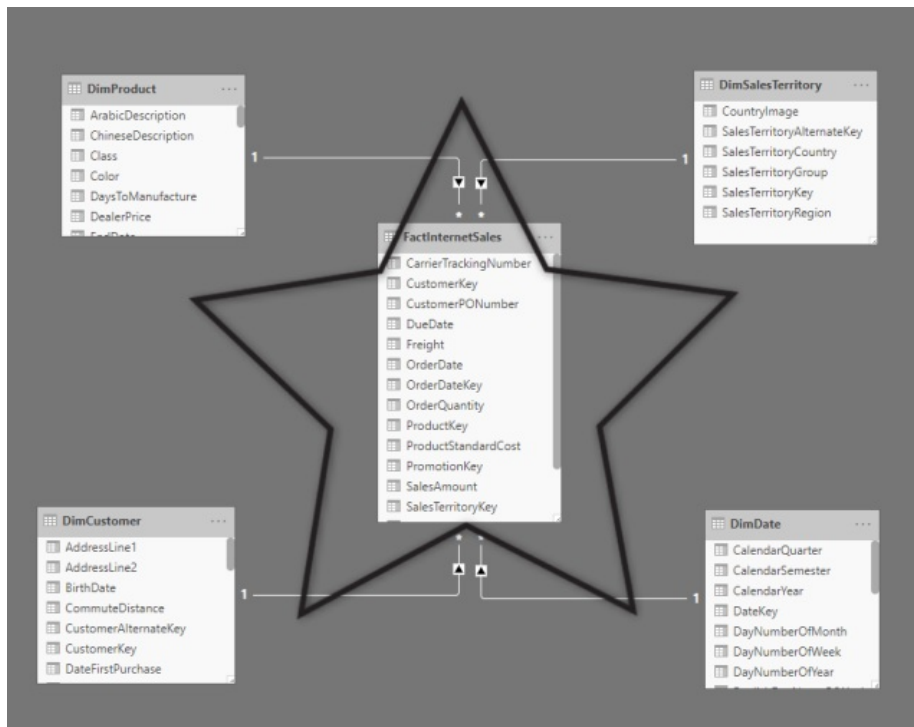


Figura 3. Ejemplo de modelo de datos en estrella en Power BI, con una tabla de hechos central y tablas de dimensiones relacionadas que facilitan la agregación y segmentación de la información. Fuente:

<https://radacad.com/power-bi-basics-of-modeling-star-schema-and-how-to-build-it/>.

En entornos como Power BI y Tableau, el modelado de datos se basa en relaciones uno a muchos o muchos a uno entre tablas de hechos y dimensiones. Este enfoque permite agregar métricas de forma flexible y segmentar resultados según atributos categóricos. Una práctica recomendable es definir un esquema estrella o copo de nieve que simplifique el mantenimiento y acelere las consultas.

Otra consideración importante es el uso de claves únicas y la identificación de campos que permitan unir tablas de manera inequívoca. La falta de integridad referencial puede producir errores de agregación y visualizaciones incompletas. Por ello, es esencial comprobar que cada tabla dispone de identificadores consistentes y que las relaciones están correctamente documentadas.

Optimización del rendimiento de carga

El rendimiento de un dashboard depende en gran medida de cómo se han preparado los datos que alimentan los gráficos. Una carga lenta o una actualización deficiente pueden afectar negativamente a la experiencia del usuario y reducir la confianza en la herramienta. Por este motivo, la optimización del rendimiento es una etapa clave del preprocesamiento.

Entre las estrategias más habituales se encuentra la reducción de la granularidad de los datos cuando no es necesaria información a nivel de detalle. Por ejemplo, en vez de cargar transacciones individuales, se pueden precalcular agregados por mes o categoría. Asimismo, conviene eliminar campos que no se utilizarán en los análisis, ya que cada columna adicional incrementa el peso de la memoria.

Otra buena práctica consiste en crear índices y particiones en las bases de datos origen cuando se trabaja con volúmenes significativos. Estas técnicas permiten acelerar la recuperación de datos y reducir el tiempo de respuesta de las consultas. Finalmente, en herramientas como Power BI, el uso de modelos importados en lugar de conexiones en directo contribuye a mejorar la fluidez del dashboard.

Consideraciones de actualización automática

La actualización periódica de los datos es un requisito imprescindible en entornos empresariales donde los indicadores se monitorizan en tiempo real, o con frecuencia diaria. Configurar estos procesos correctamente garantiza que los usuarios consultan siempre la versión más actualizada de la información y que se evitan errores derivados de desfases temporales.

Las herramientas de BI ofrecen distintas opciones para automatizar la actualización: desde la programación de actualizaciones en el servicio de Power BI hasta la publicación programada de extractos en Tableau Server. En proyectos que integran APIs o fuentes dinámicas, es habitual definir procesos ETL que descargan y transforman los datos de forma desatendida antes de incorporarlos al modelo.

Además de la programación técnica, es importante establecer alertas y mecanismos de control que verifiquen que las actualizaciones se han realizado correctamente. El registro de incidencias, la monitorización de tiempos de carga y la documentación de cada proceso de refresco forman parte de las buenas prácticas que aseguran la fiabilidad de los *dashboards* en el tiempo.

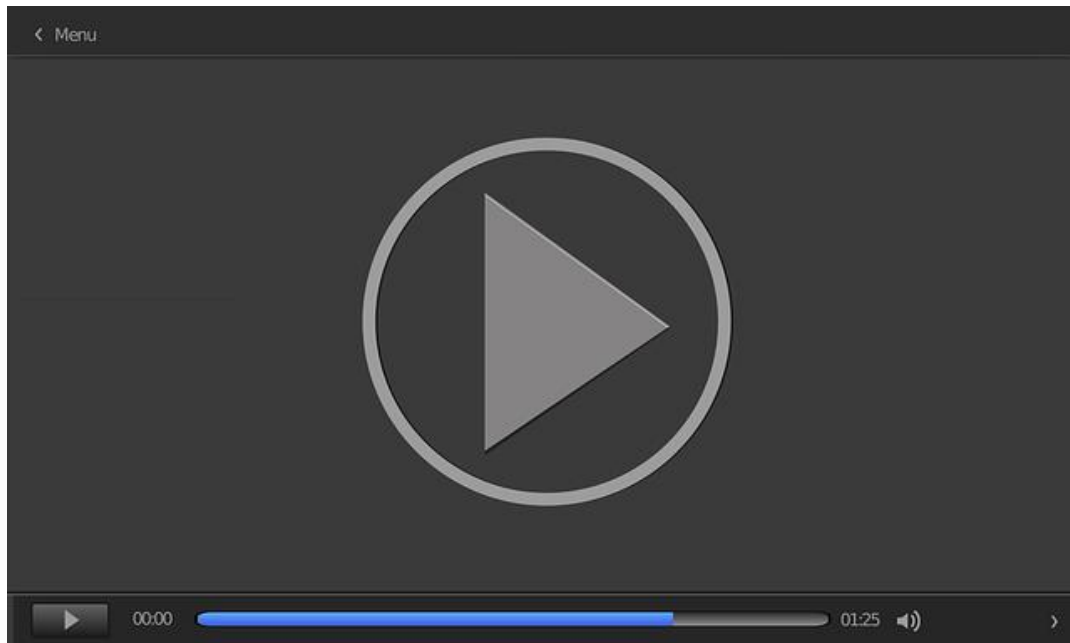
3.7. Resumen y conclusiones

La preparación de datos es mucho más que una tarea preliminar: constituye el núcleo de cualquier proyecto de visualización avanzada. A lo largo de este tema se han revisado los principios que guían la organización, limpieza, transformación, e integración de información, destacando la importancia de aplicar metodologías sistemáticas y herramientas especializadas que garanticen la calidad y consistencia del proceso.

Se han explorado las distintas opciones disponibles para realizar estas tareas, desde entornos visuales como Power BI, Tableau Prep, y OpenRefine hasta flujos de programación con R y Python. Cada herramienta aporta ventajas específicas que pueden combinarse para adaptarse a la naturaleza del proyecto y al perfil del equipo responsable. La capacidad de elegir con criterio y documentar cada transformación es una competencia esencial para asegurar la reproducibilidad y la confianza en los resultados.

Por último, se ha puesto de relieve que el éxito de un *dashboard* no depende únicamente de su diseño final, sino de la solidez de los datos que lo alimentan. Modelar tablas, optimizar el rendimiento de carga y establecer procesos de actualización automática son pasos imprescindibles para construir soluciones visuales fiables y sostenibles en el tiempo. La profesionalidad en la preparación de datos marca la diferencia entre una visualización que informa y otra que transforma la toma de decisiones.

Para ahondar más acerca de este tema, puedes revisar el vídeo *Limpieza y estructuración de datos para visualización*.



Accede al vídeo:

<https://unir.cloud.panopto.eu/Panopto/Pages/Embed.aspx?id=d276e55f-2f79-4482-8d81-b329007d95b6>

3.8. Referencias bibliográficas

Deckler, G. (2022). *Learn Power BI: A comprehensive, step-by-step guide for beginners to learn real-world business intelligence*. Packt Publishing Ltd.

Ham, Y. (2023). *Hands-On Data Cleaning and Preprocessing in Python*. Packt Publishing.

McKinney, W. (2022). *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and Jupyter*. O'Reilly Media.

OpenRefine. (2023). *OpenRefine Documentation*. <https://openrefine.org/docs/>

Tableau Software. (2023). *Tableau Prep Help*. https://help.tableau.com/current/prep/en-us/prep_get_started.htm

Wickham, H., Cetinkaya-Rundel, M., y Grolemund, G. (2023). *R for Data Science: Import, Tidy, Transform, Visualize, and Model Data* (2nd ed.). O'Reilly Media. <https://r4ds.hadley.nz/>

The Art of Data Cleaning: Transforming Messy Data into Structured Insights

Halder, N. (2023). *The art of data cleaning: Transforming messy data into structured insights*. Medium. <https://medium.com/%40HalderNilimesh/the-art-of-data-cleaning-transforming-messy-data-into-structured-insights-b20fece2c670>

Este artículo ofrece una visión clara y práctica del proceso de limpieza de datos, con énfasis en casos reales y pasos esenciales para convertir datos "sucios" en conjuntos listos para análisis y visualización.

Guide To Data Cleaning: Definition, Benefits, Components, And How To Clean Your Data

Tableau Software. (2025). *Guide to data cleaning: Definition, benefits, components, and how to clean your data*. Tableau Blog. <https://www.tableau.com/learn/articles/what-is-data-cleaning>

Guía oficial de Tableau que describe las fases del proceso de limpieza, sus beneficios y componentes, con ejemplos prácticos que refuerzan conceptos tratados en Power BI y Tableau Prep.

OpenRefine Tutorial – University of Washington

University of Washington. (2025). *OpenRefine Tutorial*. Open Data Literacy. https://odl.ischool.uw.edu/openrefine_tutorial/

Tutorial especializado que enseña cómo usar OpenRefine para limpiar metadatos desde APIs públicas, incluyendo cambio de formato, duplicados y análisis exploratorio.

Tutorial de OpenRefine: limpieza y transformación de datos abiertos

Junta de Andalucía. (s. f.). *Tutorial de OpenRefine: limpieza y transformación de datos abiertos*. Portal de Datos Abiertos de Andalucía. Disponible en <https://www.juntadeandalucia.es/datosabiertos/portal/tutoriales/usar-openrefine.html>

Guía práctica en español que muestra paso a paso cómo usar OpenRefine para explorar, limpiar y transformar datos abiertos. Incluye ejemplos aplicados y capturas que ayudan a familiarizarse con la interfaz y las funciones principales.

1. ¿Qué principio fundamental establece que cada variable debe estar en una columna y cada observación en una fila?
 - A. Modelo estrella
 - B. Tidy data
 - C. Esquema copo de nieve
 - D. Normalización de segundo nivel

2. ¿Qué herramienta permite la limpieza de datos mediante agrupamiento por similitud y transformación con expresiones regulares?
 - A. Tableau Prep
 - B. Power Query
 - C. Talend Open Studio
 - D. OpenRefine

3. ¿Cuál es una ventaja de usar Power Query en Power BI?
 - A. Permite documentar cada paso de transformación en un flujo visual
 - B. Solo se puede usar en versiones de escritorio
 - C. No admite conexiones con bases de datos
 - D. Requiere programación en Python

4. ¿Qué librería de Python es la más utilizada para manipulación de datos tabulares?
 - A. NumPy
 - B. Matplotlib
 - C. Pandas
 - D. Scikit-learn

5. ¿Cuál de las siguientes no es una práctica recomendada en la integración de datos?
- A. Verificar claves de relación
 - B. Homogeneizar formatos de fechas
 - C. Mezclar categorías distintas sin documentar
 - D. Documentar fuentes y criterios de integración
6. ¿Qué técnica mejora el rendimiento de carga en dashboards?
- A. Reducir granularidad de datos innecesaria
 - B. Incluir todas las columnas de origen
 - C. Evitar particiones e índices
 - D. Conectar siempre en modo directo
7. ¿Qué librería de R es clave para transformar datos de forma declarativa?
- A. Ggplot2
 - B. Dplyr
 - C. Shiny
 - D. Readxl
8. ¿Cuál es una razón para almacenar descargas de APIs en repositorios intermedios?
- A. Reducir espacio de almacenamiento
 - B. Evitar toda validación posterior
 - C. Garantizar la trazabilidad y disponer de históricos
 - D. Eliminar datos antiguos de inmediato

9. ¿Qué elemento es esencial al definir relaciones entre tablas?
- A. Usar nombres de campos arbitrarios
 - B. Evitar claves primarias
 - C. Asegurar unicidad de identificadores
 - D. Mezclar diferentes unidades sin conversión
10. ¿Qué funcionalidad permite programar actualizaciones automáticas de datos en Tableau Server?
- A. Tableau Prep Viewer
 - B. Tableau Prep Conductor
 - C. Tableau Public
 - D. Tableau Bridge