



ESCUELA POLITÉCNICA DE INGENIERÍA DE GIJÓN

GRADO EN INGENIERÍA INFORMÁTICA EN TECNOLOGÍAS DE LA INFORMACIÓN

Lenguajes y Sistemas Informáticos

TRABAJO FIN DE GRADO/MÁSTER Nº ???

Explotación, integración y visualización de múltiples fuentes de datos mediante un Data Lake

Mier Montoto, Juan Francisco

TUTORES:

D. Augusto Alonso, Cristian

D. Morán Barbón, Jesús

D. Vázquez Faes, Eduardo

FECHA: junio 2024

Índice de contenido

Índice de contenido	1
Índice de figuras	2
1. Introducción	3
1.1. Antecedentes	3
1.2. Motivación	3
1.3. Finalidad del proyecto	4
1.4. La empresa	4
2. Fundamento teórico	6
2.1. Definición de Data Lake	6
2.2. Procesos ETL	6
2.3. Dashboards	6
2.4. Componentes del stack	7
2.4.1. Tópico	7
2.4.2. Productor	7
2.4.3. Consumidor	7
2.4.4. Broker	7
2.4.5. Zookeeper	8
3. Planificación	9
3.1. Planificación inicial	9
3.2. Presupuesto inicial	9
4. Análisis	10
4.1. Definición del sistema	10
4.2. Análisis de alternativas	10
5. Diseño del sistema	11
5.1. Arquitectura del sistema	11
5.2. Modelo de datos	11
6. Implementación	12
7. Resultados	13
8. Conclusiones y trabajo futuro	14
Bibliografía	15

Índice de figuras

1. Introducción

1.1. Antecedentes

En la actualidad, la cantidad de datos que se generan y almacenan es cada vez mayor ¹, una tendencia que por supuesto se traduce a las empresas. Estos datos provienen de múltiples fuentes y en múltiples formatos, lo que dificulta su análisis y explotación. A esta característica de la información se le conoce como *heterogeneidad*².

A su vez, el progreso tecnológico ha permitido la creación de nuevas herramientas y técnicas que facilitan la recogida, almacenamiento y análisis de estos datos. Una de estas técnicas son los *data lakes* (ver 2.1 *Definición de Data Lake*), que permiten almacenar grandes cantidades de datos de diferentes tipos y formatos, para poder analizarlos y explotarlos de forma más eficiente.

De manera paralela, el progreso en la *ciencia de datos* ha supuesto una mejora en la visualización de estos, lo que permite un mejor análisis. La visualización de datos suele significar la existencia de uno o varios *dashboards* [1], que permiten al usuario final la consulta rápida y sencilla de los datos sin necesidad de conocimientos técnicos del sistema.

1.2. Motivación

El proyecto surge de la necesidad de la empresa (ver 1.4 *La empresa*) de recoger y analizar datos heterogéneos de todas las fuentes de las que se disponen, tanto internas (e.g. bases de datos, archivos de registros, APIs, entre otros), como externas (e.g. APIs o datos de webs de terceros, datos de fuentes públicas...).

En la actualidad, la empresa dispone de una gran cantidad de datos que se encuentran en diferentes formatos y en diferentes ubicaciones, lo que dificulta su análisis y explotación. Por otra parte, se depende de la consulta manual o de servicios de terceros (como dashboards en NewRelic o AWS CloudWatch) para poder analizar estos datos, lo que supone un coste adicional.

Además del uso interno, la empresa también quiere ofrecer a sus clientes la posibilidad de consultar estos datos de forma visual y sencilla, para que puedan analizarlos y explotarlos de

¹<https://www.statista.com/statistics/871513/worldwide-data-created/>

²<https://www.sciencedirect.com/topics/computer-science/data-heterogeneity>

forma autónoma, lo que supondría un valor añadido para los mismos. Este tipo de dashboards son diferentes a los dashboards de monitorización antes mencionados, ya que permiten al usuario final la consulta de datos de negocio, y no de infraestructura.

1.3. Finalidad del proyecto

El objetivo de este sistema es centralizar y unificar las fuentes de datos heterogéneas cuya consulta se realiza de manera manual, con la finalidad de analizar los datos de forma más eficiente.

El cumplimiento de este objetivo permitirá a la empresa obtener una serie de beneficios:

- una eliminación del tiempo invertido en la consulta manual de los datos.
- una reducción de los costes de las plataformas de terceros.
- una mejora en la toma de decisiones, al poder analizar los datos de forma más eficiente.
- una mejora en la calidad de los servicios ofrecidos a los clientes, al poder ofrecerles la posibilidad de consultar los datos de forma visual y sencilla.
- el beneficio económico que supondría la venta de este servicio a los clientes.

Además de la mejora de los procesos ya existentes, la explotación mediante esta herramienta abrirá la puerta a nuevas posibilidades de análisis y explotación de los datos, como la detección de anomalías en la infraestructura o la predicción de patrones y eventos futuros.

1.4. La empresa

Okticket es una startup nacida en Gijón en 2017 cuyo producto principal es un servicio software que reduce los costes y el tiempo que invierten las empresas en contabilizar y manejar los gastos de viaje de los profesionales mediante el escaneo automático de tickets y notas de gastos.

Sus oficinas principales (incluyendo la zona de desarrollo) se encuentran en el Parque Tecnológico de Gijón, aunque cuenta con un número de sedes creciente en varios países, como Francia, Portugal o, más recientemente, México.

Okticket es una de las empresas que más crecen tanto del sector como del propio Parque Tecnológico. Debido a este rápido crecimiento, el equipo está en constante desarrollo y cambio, tanto aquí en España como en el resto de sedes. Este crecimiento se refleja en la recepción de un gran número de galardones y reconocimientos ^{3 4 5 6}

La parte principal del negocio es el núcleo del software como servicio (Software as a Service en inglés, en adelante *SaaS*), es decir, la aplicación completa tanto para administradores como para empleados. Este SaaS se oferta a empresas de cualquier tamaño, cuyo precio final varía en función del número de usuarios, las características e integraciones que requiera la empresa cliente y el soporte que se ofrezca.

Recientemente se han añadido nuevas propuestas a la cartera de servicios ofertada por Okticket, como la OKTCard - una tarjeta inteligente que gestiona automáticamente los gastos, así como la inclusión de nuevos “módulos” de gestión de gastos y viajes.

³ Okticket en el especial startups 2023 de Forbes (LinkedIn)

⁴ Arcelor y Okticket, premios nacional de Ingeniería Informática (EL COMERCIO)

⁵ Okticket recibe el sello Pyme Innovadora (okticket.es)

⁶ Okticket, empresa emergente certificada (okticket.es)

2. Fundamento teórico

2.1. Definición de Data Lake

Los *data lakes*¹ son almacenes de datos que almacenan grandes cantidades de datos de manera no estructurada [1]. A diferencia de los *data warehouses*, los *data lakes* no tienen un esquema definido, lo que permite almacenar datos de cualquier tipo y formato. Esto permite almacenar grandes cantidades de datos sin tener que definir un esquema de antemano, lo que puede ser útil en algunos casos. Sin embargo, esto también puede ser un inconveniente, ya que no se puede realizar un análisis de negocio de los datos almacenados en un data lake sin antes transformarlos en un formato estructurado.

2.2. Procesos ETL

Los procesos ETL [1] son procesos que combinan datos de múltiples fuentes en un único destino, transformando los datos en un formato común. Estos procesos se utilizan para extraer datos de diferentes fuentes, transformarlos en un formato común y cargarlos en un dashboard para que se muestren de manera visual.

Los procesos ETL deben adaptarse a la estructura de los datos de la fuente de origen, ya que cada sistema procesa diferentes tipos de datos, de diferentes fuentes y en diferentes formatos.

Una de las características clave de los procesos ETL es que sean escalables, ya que los datos que se muestran en los dashboards suelen ser datos que se generan de manera continua, y por lo tanto los procesos ETL deben ser capaces de procesar grandes cantidades de datos de manera eficiente. En ocasiones, los procesos ETL se pueden realizar en *streaming*, lo que significa que los datos se procesan en tiempo real a medida que se generan.

2.3. Dashboards

La palabra *dashboard*, que traducido de manera literal significa *cuadro de mandos*, es un término que se utiliza para referirse a cualquier interfaz gráfica que muestre información relevante de manera visual sobre un proceso o negocio. Aunque el término se utiliza en

¹<https://aws.amazon.com/es/what-is/data-lake/>

muchos ámbitos (puede incluir indicadores comerciales, de producción, de marketing, de calidad, de recursos humanos...)

En el ámbito de este proyecto, los dashboards reflejan en tiempo real el rendimiento de actividades o procesos de negocio, y se utilizan para tomar decisiones informadas sobre los mismos. En el caso de una empresa, un dashboard puede mostrar desde el rendimiento de la plataforma en tiempo real hasta un reflejo del de las ventas, y permitir a los directivos tomar decisiones informadas sobre el futuro de la empresa.

2.4. Componentes del stack

En el *stack* tecnológico escogido para el proyecto se manejan diferentes términos y conceptos que son necesarios desarrollar para entender el funcionamiento del sistema.

2.4.1. Tópico

Un tópico es una categoría a la que se envían los mensajes. Al proceso de enviar mensajes a un tópico se le llama *publicar*, y al proceso de leer mensajes de un tópico se le llama *consumir*.

2.4.2. Productor

El productor es el componente responsable de crear y enviar mensajes al cluster de Kafka. Está separado del resto de los componentes y produce mensajes de manera asíncrona y rápida.

2.4.3. Consumidor

El consumidor es el componente responsable de leer los mensajes producidos por el productor. Está suscrito a un Tópico a través del broker y consume los mensajes.

2.4.4. Broker

El broker es el componente responsable de recibir los mensajes producidos por el productor y enviarlos a los consumidores. Es el intermediario entre los productores y los consumidores.

2.4.5. Zookeeper

Zookeeper es un servicio de coordinación distribuida que se utiliza para gestionar y coordinar los brokers de Kafka. Se encarga de mantener la información de los brokers y de los tópicos.

3. Planificación

3.1. Planificación inicial

3.2. Presupuesto inicial

4. Análisis

4.1. Definición del sistema

4.2. Análisis de alternativas

5. Diseño del sistema

5.1. Arquitectura del sistema

5.2. Modelo de datos

6. Implementación

7. Resultados

8. Conclusiones y trabajo futuro

Bibliografía

- [1] J. Mier, “Presentación de datos: dashboards y procesos ETL.” Primera entrega de teoría de la asignatura Inteligencia de Negocio, EPI Gijón, curso 23-24, 2023.
- [2] J. Mier, “latexTemplate.” <https://github.com/miermontoto/latexTemplate>, 2024. Plantilla de L^AT_EX personal para trabajos académicos.
- [3] J. Mier, “Minería de anomalías.” Segunda entrega de teoría de la asignatura Inteligencia de Negocio, EPI Gijón, curso 23-24, 2024.