


Bases de datos para datos masivos

Contenido de la materia

Tema 1





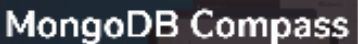



MÉTODOS DE CAPTURA DE INFORMACIÓN		
Origen y calidad de los datos	Organización de los datos	Casos de estudio
<p>Evaluación de calidad:</p> <ul style="list-style-type: none">- Compleitud: grado en el que los valores se encuentran en un conjunto de datos.- Credibilidad: nivel de fiabilidad del organismo que proporciona el conjunto de datos.- Consistencia: grado en el que los datos carecen de contradicciones.- Interpretabilidad: grado en el que los datos deben ser interpretados por una persona.- Precisión: nivel de exactitud del valor.	<p>Ficheros planos:</p> <ul style="list-style-type: none">- CSV (<i>comma separated values</i>): RFC 4180. Define un registro por línea y separa los campos por comas.- JSON: RFC 7159 y ECMA-404. Describe objetos encapsulados por llaves y listas por corchetes.- XML: descrito en el estándar XML 1.0 de W3C. Permite almacenar información de forma legible utilizando etiquetas.	<p>Procesamiento de sitio web sobre cursos online.</p>
<p>Niveles de abstracción:</p> <ul style="list-style-type: none">- Datos: conjunto de hechos discretos y objetivos sobre un evento.- Información: datos con significado.- Conocimiento: combinación de información contextualizada, experiencias, valores e intuición.	<p>Bases de datos:</p> <ul style="list-style-type: none">- Conjunto de datos persistentes, utilizados por sistemas de aplicación.- En el modelo Entidad-Relación (E/R) una entidad es cualquier objeto repensado en la BBDD y un vínculo representa relaciones entre ellos.- La unidad básica de almacenamiento es el campo, agrupados en registros y estos en ficheros almacenados.	<p>Procesamiento de logs de servidor web.</p>
<p>Fuentes de datos:</p> <ul style="list-style-type: none">- Captura manual: encuestas y observaciones.- Análisis de documentos estructurados: estructurados (HTML) y sin formato (lenguaje natural).- Salida de aplicaciones: logs o bases de datos.- Sensores: dispositivos de medición.- Datos de acceso público: gubernamentales y servicios web públicos.	<p>Bases de datos relacionales y SQL:</p> <ul style="list-style-type: none">- Los ficheros almacenados se representan en forma de tablas (relaciones), con columnas (campos) y filas (registros).- El estándar SQL define un lenguaje para la consulta y modificación de los datos.- El comando SELECT permite consultar información de tablas.- Los comandos INSERT, UPDATE y DELETE permiten la inserción, edición y eliminación de registros respectivamente.	<p>API de acceso a transacciones bancarias.</p>
		<p>Almacenamiento de información sobre productos en un fichero CSV.</p>
		<p>Representación de información geolocalizada en formato JSON.</p>
		<p>Almacenamiento de información sobre clientes de una base de datos relacional.</p>
		

Tema 2

NoSQL		
Definición	NoSQL vs. SQL	Bases de datos NoSQL
<p>El movimiento NoSQL incluye todas las bases de datos con arquitectura distinta a la utilizada en sistemas relacionales tradicionales.</p>	<p>NoSQL:</p> <ul style="list-style-type: none">- Es una tecnología que se lleva utilizando desde los años 60, aunque el nombre fue acuñado en 2009.- Las bases de datos NoSQL se caracterizan por su escalabilidad horizontal y sencillez, evitan cuellos de botella, permiten manejar grandes volúmenes de datos y su puesta en funcionamiento es más asequible.	<p>Cada base de datos siguiente corresponde a un tipo de NoSQL</p>
<p>Tipos</p> <p>Las bases de datos NoSQL pueden categorizarse en cuatro grupos:</p> <ul style="list-style-type: none">- Clave-valor simples.- Clave-valor sofisticadas.- Basadas en documentos.- Basadas en grafos.	<p>Desventajas:</p> <ul style="list-style-type: none">- La falta de madurez, en algunos casos, es un problema porque lleva consigo escasa documentación y comunidades de usuarios muy reducidas.- Problemas de compatibilidad con otras herramientas, aunque esto es un problema que tiende a desaparecer debido a la apuesta que hacen las grandes compañías al integrar este tipo de productos en sus servicios.	<p>Apache Cassandra</p> <p>Es una base de datos distribuida de código abierto escrita en Java. Todos sus nodos actúan por igual, agrupándose en anillos. Permite sistemas de réplicas y el acceso a los datos se hace a través de CQL.</p>
<p>Teorema CAP</p> <p>También llamado teorema de Brewer, señala que todo sistema distribuido no puede garantizar a la vez que haya consistencia, disponibilidad y tolerancia a particiones (<i>consistency-availability-partition tolerance</i>).</p>	<p>Patrones de diseño:</p> <p>MongoDB es un buen ejemplo de base de datos NoSQL que aplica relaciones entre documentos para aprovechar la utilidad de los modelos relacionales, sin llegar a ser un relacional en toda regla.</p> <p>Relaciones que se crean en MongoDB:</p> <ul style="list-style-type: none">- Relaciones entre documentos uno a uno.- Relación uno a muchos con documentos embebidos.- Relación uno a muchos con documentos referidos.	<p>Neo4J</p> <p>Es una base de datos basada en grafos, compatible con ACID. Es accesible desde <i>software</i> escrito en otros lenguajes usando Cypher Query Language, a través de un punto HTTP transaccional.</p>
		<p>MongoDB</p> <p>Es una base de datos basada en documentos, la cual almacena los datos en esta estructura. El conjunto de documentos se denomina colecciones y el conjunto de estos conforman la base de datos.</p> <p>El formato de almacenamiento es BSON.</p>

Se utiliza para tipificar las bases de datos NoSQL


Tema 3


TRATAMIENTO DE DATOS EN MONGODB				
Software de apoyo	Flexibilidad del modelo	Inserción de datos	Lecturas y consultas	Actualización de datos
Existen muchas aplicaciones, tanto internas como externas, que facilitan el uso del MongoDB.	El esquema de MongoDB es flexible permitiendo modelos de datos totalmente diferentes.	El comando <code>db.<collection>.insert</code> recibe el objeto a añadir en formato JSON. 	El comando principal para la consulta de documentos es <code>db.<collection>.find</code> . 	El comando <code>db.<collection>.save</code> actualiza un documento si ya existe en la colección. 
MongoBooster es una herramienta GUI multiplataformas que permite la construcción de consultas fluidas y que posee un gran número de herramientas muy útiles.  MongoBooster	Dos documentos de una colección pueden tener atributos totalmente diferentes.	El comando <code>db.<collection>.save</code> crea un documento en la <i>collection</i> si este no existía previamente.	Find recibe dos parámetros: el criterio de búsqueda y la proyección de atributos a mostrar.	El comando <code>db.<collection>.update</code> permite cambiar atributos específicos de un conjunto de documentos.
MongoDB Compass es la solución propietaria de este tipo de herramienta. 	En la práctica, la mayoría de documentos de una colección comparten una estructura similar, aunque esto no es una regla fija.		El criterio de búsqueda puede incluir operaciones como menor que (<code>\$lt</code>) y mayor que (<code>\$gt</code>), y en (<code>\$in</code>).	El comando <code>update</code> recibe tres parámetros: criterio, acción y opciones. 
	<pre>{ "firstName": "Ray", "lastName": "Williams", "joined": { "month": "January", "day": 12, "year": 1982 } }, { "firstName": "John", "lastName": "Jones", "joined": { "month": "April", "day": 25, "year": 1986 } }</pre>		Los resultados de una consulta pueden modificarse para ser ordenados, truncados u omitidos.	El comando para eliminar documento es <code>db.<collection>.remove</code> . Recibe un parámetro con el criterio de búsqueda. 
			El comando <code>db.<collection>.findOne</code> permite obtener un solo documento. 	

Tema 4

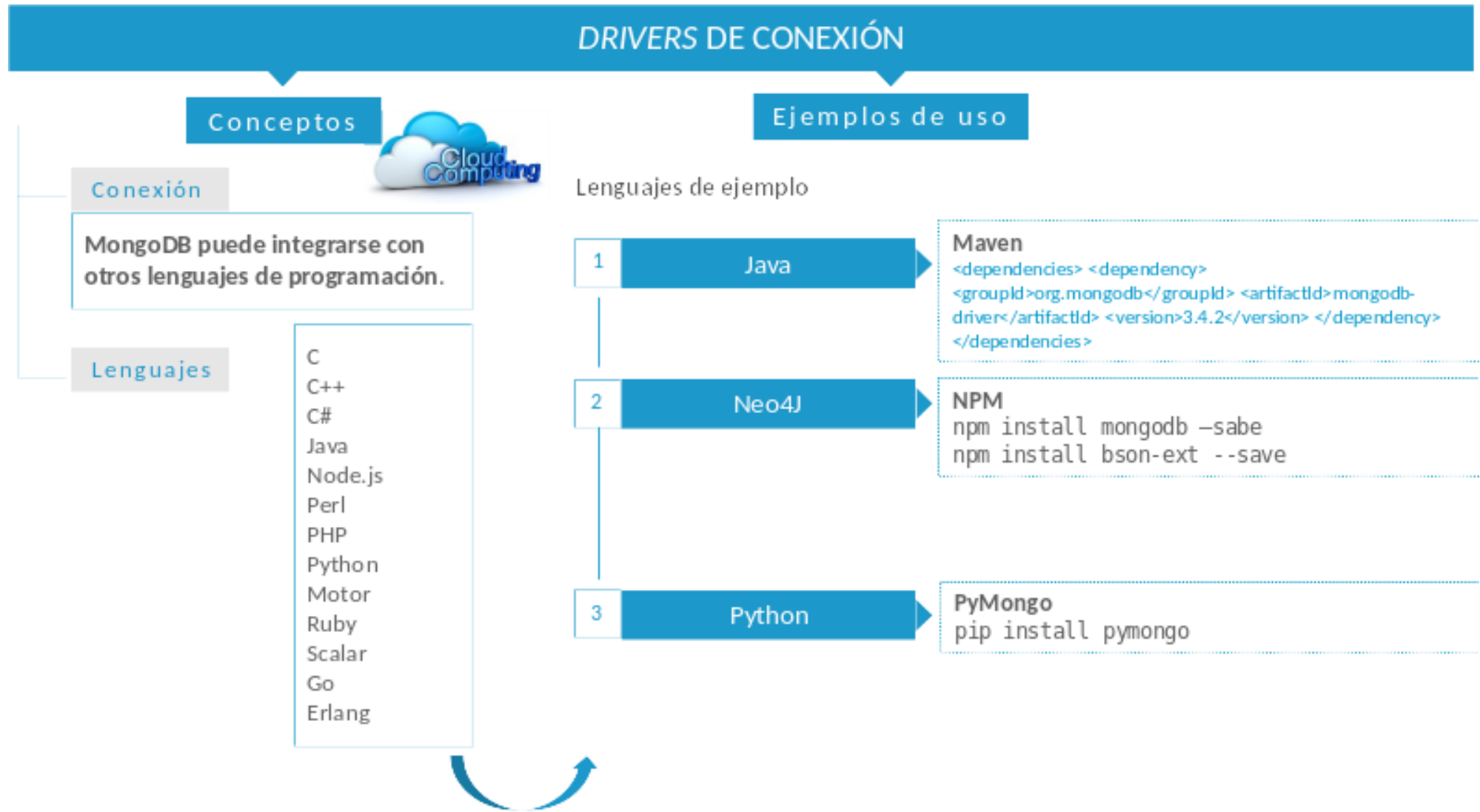
AGREGACIÓN		
Concepto	Map-Reduce	Aggregation Framework
Las funciones de agregación son útiles para agrupar datos y obtener resultados a partir de ellos.	Map-Reduce está formada por dos funciones. Una función <i>map</i> que genera los pares clave-valor y otra <i>reduce</i> , que opera sobre ellos.	El <i>framework</i> de agregación está modelado en varias etapas que transforman los documentos en un resultado agregado.
En la base de datos relacionales, este tipo de consultas se realizan con el operador group by (para agrupar) y sum , count , etc., para realizar cualquier operación con sus datos.	El lenguaje de programación utilizado es JavaScript.	Las etapas mas usadas son la \$match , encargada de la búsqueda, y \$group , la que agrupa y hace los cálculos.
Existen dos métodos de agregación en MongoDB: Map-Reduce y Aggregate .	Los parámetros de agregación en la función Map-Reduce se especifican como un objeto JSOM.	Es obligatorio especificar el campo _id en el nuevo documento, ya que estos son los valores a agrupar.
Existen operaciones específicas de agregación: count , distinct y group .	La mayor ventaja de la utilización de este método de agregación es la potencia que ofrece, ya que se pueden personalizar totalmente cada una de sus funciones.	Además de los operadores de búsqueda, tiene muchos otros que hacen la vida mas fácil al trabajar, <i>array</i> , <i>string</i> , etc.

Tema 5

GESTIÓN DE MONGODB			
Seguridad	Respaldo	Rendimiento	Sharding
<p>Autenticación y listas de acceso: permiten controlar permisos de interacción con bases de datos.</p>	<p>La herramienta Mongodump se utiliza para crear una copia de respaldo de una base de datos en MongoDB.</p>	<p>Los índices permiten ejecutar consultas de forma eficiente.</p>	<p>Sharding: método utilizado por MongoDB para distribuir datos en varios servidores y aumentar así su capacidad.</p>
<p>Para restringir conexiones desde una interfaz en específico se debe indicar la dirección IP en el parámetro <code>-Bind_ip</code></p>	<p>La herramienta Montorestore permite restaurar una copia de respaldo en el sistema MongoDB.</p>	<div></div>	<p>Cada <i>shard</i> está encargado de almacenar la información. Cada <i>shard</i> es un replica set, por lo que brinda alta disponibilidad y consistencia en datos.</p>
<p>La seguridad del entorno MongoDB puede incrementarse con el uso de un <i>firewall</i>.</p>	<p>MongoDB brinda capacidades de replicación, para lograr reducirse e incrementar la capacidad.</p>	<p>Este aumento de rendimiento se ve reflejado en el número de documentos analizados y, por consiguiente, en el tiempo de respuesta.</p>	



Tema 6



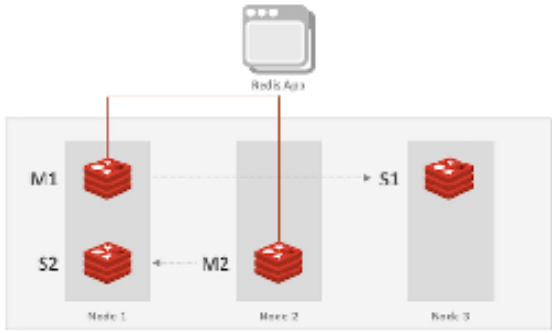


Tema 7

Apache Cassandra		
Definición	Generalidades	CQL3
<p>Cassandra es un producto de Apache y su sistema de almacenamiento o motor de base de datos es de código abierto, distribuido y descentralizado. Su fin principal es gestionar grandes volúmenes de datos estructurados repartidos por distintos servidores.</p>	<p>Arquitectura:</p> <ul style="list-style-type: none">- Cassandra está diseñada para trabajar en entornos <i>big data</i> sobre múltiples nodos evitando puntos de fallo.- Dispone de un sistema distribuido de igual a igual en todos sus nodos donde los datos se distribuyen entre todo el clúster.- Los nodos de un clúster pueden aceptar solicitudes de lectura y escritura, independientemente de dónde se encuentren realmente los datos en dicho clúster.- Si un nodo falla, las solicitudes de lectura/escritura pasan a ser atendidas por otros nodos de la red.	<p>Keyspaces, tablas y columnas</p> <p>Manipulación de datos</p> <ul style="list-style-type: none">• Tipos de datos.• Contadores.• Uso de <i>timestamp</i>.• Uso de <i>date</i>.• Definir duraciones.• Colecciones: <i>maps</i>, <i>set</i> y <i>list</i>.• Tuplas.• UDT.• Operaciones de INSERT, UPDATE, DELETE.
<p>Base de datos NoSQL de tipo clave-valor</p>	<p>Modelo de datos:</p> <ul style="list-style-type: none">- Los patrones de acceso a los datos y las consultas de la aplicación determinan la estructura y organización de los datos que luego se utilizan para diseñar las tablas de la base de datos.- Cada consulta está respaldada por una tabla y, por ello, los datos se duplican en varias tablas en un proceso conocido como desnormalización.- La duplicación de datos y un alto rendimiento de escritura se utilizan para lograr un alto rendimiento de lectura.	<p>Otras funcionalidades</p> <ul style="list-style-type: none">• TTL.• WriteTime.• Índices.
<p>Cómo aplica el teorema CAP</p>		<p>Instalación sobre Linux y Windows</p>
<p>Conocer los principales componentes de Cassandra y describir el sistema de replicación, la lógica de las operaciones de escritura y las de lectura.</p>		

Tema 8

Neo4j		
Definición	Principales componentes	CQL (Cypher Query Language)
<p>Neo4j es una base de datos de grafos cuyo lenguaje de consultas es CQL (Cypher Query Language). Neo4j está escrito en Java.</p>	<p>Nodo:</p> <ul style="list-style-type: none"> - Las entidades reciben el nombre de nodos y cada uno de ellos es muy similar a una instancia de un objeto (es decir, tienen propiedades). 	<p>Esquemas, grafos y Cypher Query Language</p>
<p>Base de datos NoSQL de tipo grafo</p>	<p>Relaciones:</p> <ul style="list-style-type: none"> - Las relaciones, por su parte, se conocen como vértices y también tienen propiedades, siendo la dirección o sentido el más importante de ellos. Una relación conecta dos nodos y los organizan en estructuras. 	
<p>Cómo aplica el teorema CAP</p>	<p>Propiedades:</p> <ul style="list-style-type: none"> - Nodo y relaciones permiten organizar los datos de tal manera que sea posible encontrar patrones de información existente entre dichos nodos. - Las propiedades son pares de nombre-valor utilizados para agregar cualidades a los nodos y las relaciones. 	<p>Manipulación del grafo</p> <ul style="list-style-type: none"> • Tipos de datos y operadores. • Nodos. • Relaciones. • Propiedades. • Modelado. • Cláusulas CREATE, MATCH, INSERT, DELETE, REMOVE, WHERE, SKIP, LIMIT.
<div style="display: flex; align-items: center;"> <div style="text-align: center;"> <p>Nodo</p> </div> <div style="margin: 0 20px;"> <p>ACTED_IN</p> </div> <div style="text-align: center;"> <p>Nodo</p> </div> </div> <p style="text-align: center;">Relación (Vértice)</p>		

Tema 9

Redis		
Definición	Principales conceptos	Redis-cli
<p>Redis es un popular servidor de bases de datos de múltiples modelos. Redis va más allá de la base de datos NoSQL para proporcionar varias capacidades avanzadas que necesitan las aplicaciones modernas.</p>	<p>Base de datos multimodelo:</p> <ul style="list-style-type: none">- Redis proporciona una funcionalidad completa de varios modelos a través de sus módulos. El uso de Redis como una base de datos de múltiples modelos permite una mayor flexibilidad para los desarrolladores de aplicaciones dentro de una organización.	<p>Clúster, nodo, tipos de datos, estructuras y módulos</p>  <p>Redis® Cluster</p>
<p>Base de datos NoSQL de tipo clave-valor en memoria</p>	<p>Almacenamiento en memoria:</p> <ul style="list-style-type: none">- Redis mantiene los datos en la memoria para un acceso rápido y persiste los datos en el almacenamiento, así como la replicación de los contenidos en la memoria para escenarios de producción de alta disponibilidad.	
<p>Cómo aplica el teorema CAP</p>	<p>Casos de uso/estructura:</p> <ul style="list-style-type: none">- Cada estructura de datos tiene un caso de uso o escenario diferente para el que se adapta mejor.- Además de las estructuras de datos, Redis también admite el patrón Publicar/Suscribir (Pub/Sub) y patrones adicionales que hacen que sea adecuada para aplicaciones modernas con uso intensivo de datos.	<p>Estructuras de datos</p> <ul style="list-style-type: none">- <i>Strings.</i>- <i>Lists.</i>- <i>Sets.</i>- <i>Sorted sets.</i>- <i>Hashes.</i>- <i>Bit arrays.</i>- <i>Streams.</i>- <i>HyperLogLogs.</i>
<p>Los principales elementos del modelo de datos Redis son:</p> <ul style="list-style-type: none">- Tipos de datos.- Estructuras.- Módulos.		<p>Módulos</p> <ul style="list-style-type: none">- Redis Graph.- RedisSearch.- Redis TimeSeries.- RedisJSON. <p>Drivers</p> <p>Python, Java, PHP, NodeJS, C, C#.</p>
		<p>Instalación entorno de test con Docker</p>

Tema 10

