

 **IBM Bluemix** ¡Desarrolla en la nube en un click![Comience su prueba gratuita](#)

developerWorks en español Temas Técnicos Information mgmt Biblioteca técnica

# ¿Qué es Big Data?

## Todos formamos parte de ese gran crecimiento de datos

Debido al gran avance que existe día con día en las tecnologías de información han tenido que enfrentar a nuevos desafíos que les permitan analizar, descubrir lo que sus herramientas tradicionales reportan sobre su información, al mismo tiempo que los últimos años el gran crecimiento de las aplicaciones disponibles en internet (redes sociales, etc.) han sido parte importante en las decisiones de negocio de las empresas. Este artículo tiene como propósito introducir al lector en el concepto de Big Data y sus características de los componentes principales que constituyen una solución de Big Data.

Ricardo Barranco Fragoso es IT Specialist para Information Management de IBM Software Group desde diciembre de 2010. De las principales soluciones en las que está enfocado se encuentran: Master Data Management, Identity Insight y Big Data; y como este rol participa activamente apoyando las actividades de desarrollo de aplicaciones de experiencia y previo a su incorporación a IBM, Ricardo ha participado en diversos proyectos como Desarrollador Senior, Líder Técnico y Arquitecto Jr. utilizando principalmente tecnología Java para Internet Applications. Ricardo cuenta con una Licenciatura en Computación por la Universidad de Iztapalapa.

18-06-2012

---

## 1. Introducción

El primer cuestionamiento que posiblemente llegue a su mente en este momento es ¿Qué es Big Data y por qué se ha vuelto tan importante? pues bien, en términos generales podríamos referirnos como a la tendencia en el avance de la tecnología que ha abierto

las puertas hacia un nuevo enfoque de entendimiento y toma de decisiones, la cual es utilizada para describir enormes cantidades de datos (estructurados, no estructurados y semi estructurados) que tomaría demasiado tiempo y sería muy costoso cargarlos a un para su análisis. De tal manera que, el concepto de Big Data aplica información que no puede ser procesada o analizada utilizando pro tradicionales. Sin embargo, Big Data no se refiere a alguna cantidad usualmente utilizado cuando se habla en términos de petabytes y e Entonces ¿Cuánto es demasiada información de manera que sea e procesada y analizada utilizando Big Data? Analicemos primeramei

*Gigabyte* =  $10^9$  = 1,000,000,000

*Terabyte* =  $10^{12}$  = 1,000,000,000,000

*Petabyte* =  $10^{15}$  = 1,000,000,000,000,000

*Exabyte* =  $10^{18}$  = 1,000,000,000,000,000,000

Además del gran **volumen** de información, esta existe en una gran pueden ser representados de diversas maneras en todo el mundo, dispositivos móviles, audio, video, sistemas GPS, incontables senso industriales, automóviles, medidores eléctricos, veletas, anemómetros pueden medir y comunicar el posicionamiento, movimiento, vibración y hasta los cambios químicos que sufre el aire, de tal forma que las analizan estos datos requieren que la **velocidad** de respuesta sea l lograr obtener la información correcta en el momento preciso. Estas principales de una oportunidad para Big Data.

Es importante entender que las bases de datos convencionales son relevante para una solución analítica. De hecho, se vuelve mucho n en conjunto con la plataforma de Big Data. Pensemos en nuestras i

derecha, cada una ofrece fortalezas individuales para cada tarea. Como un beisbolista sabe que una de sus manos es mejor para lanzar la pelota que para atraparla; puede ser que cada mano intente hacer la actividad de la otra, pero el resultado no será el más óptimo.

## 2. ¿De dónde proviene toda esa información?

Los seres humanos estamos creando y almacenando información a una velocidad cada vez más en cantidades astronómicas. Se podría decir que si todos los datos del último año fueran guardados en CD's, se generaría una gran torre de la Luna y de regreso.

Esta contribución a la acumulación masiva de datos la podemos encontrar en industrias, las compañías mantienen grandes cantidades de datos reuniendo información acerca de sus clientes, proveedores, operaciones y de la manera que sucede con el sector público. En muchos países se administran grandes cantidades de datos que contienen datos de censo de población, registros médicos, etc. Si a todo esto le añadimos transacciones financieras realizadas en línea, datos de teléfonos móviles, análisis de redes sociales (en Twitter son cerca de 12 Tera bytes por día) y Facebook almacena alrededor de 100 Petabytes de datos geográficos mediante coordenadas GPS, en otras palabras, todas las actividades que la mayoría de nosotros realizamos varias veces al día con nuestros smartphones generan alrededor de 2.5 quintillones de bytes.

$1 \text{ quintillón} = 10^{30} = 1,000,000,000,000,000,000,000,000,000,000$

De acuerdo con un estudio realizado por Cisco[1], entre el 2011 y el 2015 el tráfico de datos móviles crecerá a una tasa anual de 78%, así como el número de dispositivos móviles conectados a Internet excederá el número de habitantes de las Naciones Unidas. Las Naciones Unidas proyectan que la población mundial alcanzará

2016 de tal modo que habrá cerca de 18.9 billones de dispositivos a escala mundial, esto conllevaría a que el tráfico global de datos mude a Exabytes mensuales o 130 Exabytes anuales. Este volumen de tráfico equivale a 33 billones de DVDs anuales o 813 cuatrillones de mensajes.

Pero no solamente somos los seres humanos quienes contribuimos con un enorme volumen de información, existe también la comunicación denominada M2M (machine-to-machine) cuyo valor en la creación de grandes cantidades de datos también es muy importante. Sensores digitales instalados en contenedores para monitorear la ruta generada durante una entrega de algún paquete y que están siendo utilizados por las compañías de transportación, sensores en medidores eléctricos que miden el consumo de energía a intervalos regulares para que sea enviada esa información a las compañías del sector energético. Se estima que hay más de 30 millones de dispositivos interconectados en distintos sectores como automotriz, transportación comercial, etc. y se espera que este número crezca en un 30% anual.

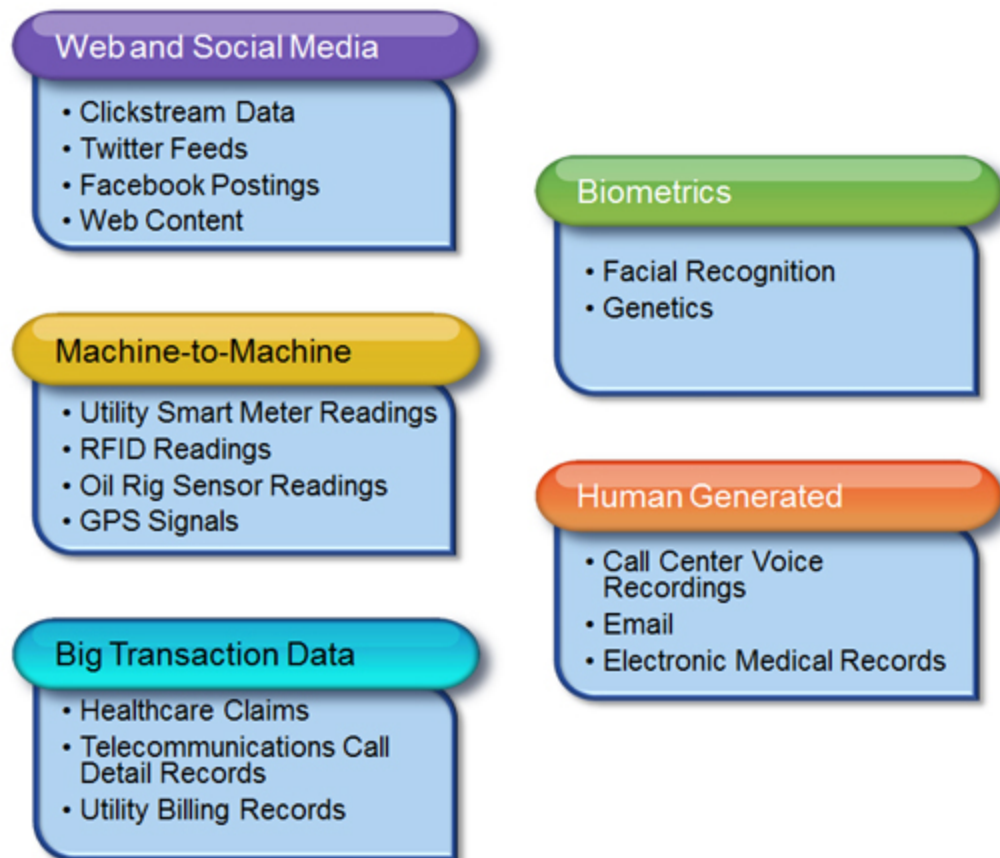
### 3. ¿Qué tipos de datos debo explorar?

Muchas organizaciones se enfrentan a la pregunta sobre ¿qué información debo analizar?, sin embargo, el cuestionamiento debería estar enfocado en el problema que se está tratando de resolver.[2]

Si bien sabemos que existe una amplia variedad de tipos de datos y que una buena clasificación nos ayudaría a entender mejor su representación, aun así, estas categorías puedan extenderse con el avance tecnológico.

#### Figura 1. Tipos de datos de Big Data[2]

## Big Data Types



1.- *Web and Social Media*: Incluye contenido web e información que sociales como Facebook, Twitter, LinkedIn, etc, blogs.

2.- *Machine-to-Machine (M2M)*: M2M se refiere a las tecnologías que otros dispositivos. M2M utiliza dispositivos como sensores o medición evento en particular (velocidad, temperatura, presión, variables químicas como la salinidad, etc.) los cuales transmiten a través de inalámbricas o híbridas a otras aplicaciones que traducen estos eventos significativa.

3.- *Big Transaction Data*: Incluye registros de facturación, en telecom detallados de las llamadas (CDR), etc. Estos datos transaccionales formatos tanto semiestructurados como no estructurados.

4.- *Biometrics*: Información biométrica en la que se incluye huellas (

retina, reconocimiento facial, genética, etc. En el área de seguridad biométricos han sido información importante para las agencias de ir

5.- *Human Generated*: Las personas generamos diversas cantidades de información que guarda un call center al establecer una llamada telefónica, correos electrónicos, documentos electrónicos, estudios médicos, e

## 4. Componentes de una plataforma Big Data

Las organizaciones han atacado esta problemática desde diferentes ángulos. Las montañas de información han generado un costo potencial al no estar asociado. Desde luego, el ángulo correcto que actualmente tiene el mayor éxito y popularidad para analizar enormes cantidades de información es la plataforma abierta *Hadoop*.

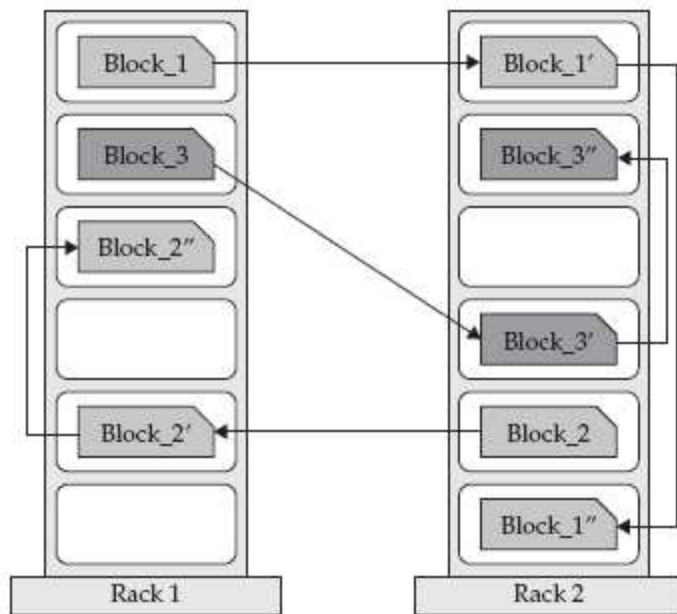
*Hadoop* está inspirado en el proyecto de Google File System (GFS) y su programación *MapReduce*, el cual consiste en dividir en dos tareas para manipular los datos distribuidos a nodos de un clúster logrando un procesamiento eficiente.[5] *Hadoop* está compuesto de tres piezas: *Hadoop* (HDFS), *Hadoop MapReduce* y *Hadoop Common*.

### ***Hadoop Distributed File System (HDFS)***

Los datos en el clúster de *Hadoop* son divididos en pequeñas piezas distribuidas a través del clúster; de esta manera, las funciones map son ejecutadas en pequeños subconjuntos y esto provee de la escalabilidad y procesamiento de grandes volúmenes.

La siguiente figura ejemplifica cómo los bloques de datos son escritos y cómo cada bloque es almacenado tres veces y al menos un bloque se almacena en un diferente rack para lograr redundancia.

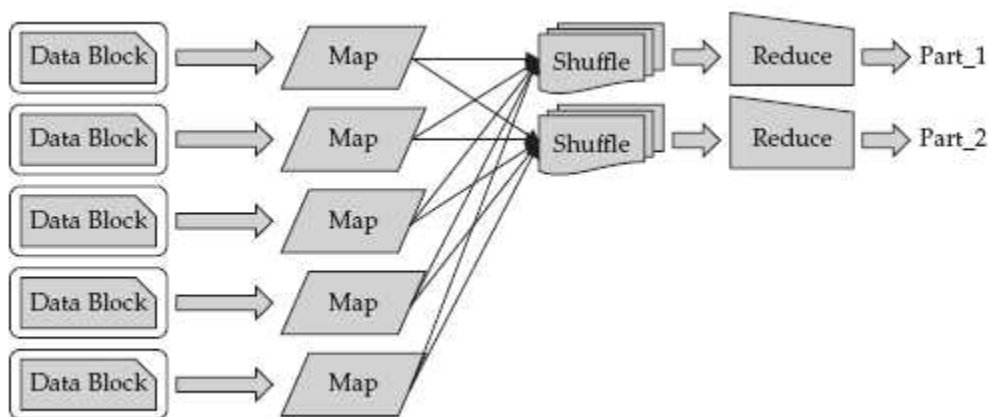
### **Figura 2. Ejemplo de HDFS**



## Hadoop MapReduce

*MapReduce* es el núcleo de Hadoop. El término MapReduce en realidad son dos procesos separados que Hadoop ejecuta. El primer proceso map, lee los datos y los convierte en otro conjunto, donde los elementos individuales son tuplas (pares de llave/valor). El proceso reduce obtiene la salida de la entrada y combina las tuplas en un conjunto más pequeño de las mapeadas. La intermedia es la denominada *Shuffle* la cual obtiene las tuplas del mapeo que cada nodo procesará estos datos dirigiendo la salida a una tarea reduce.

La siguiente figura ejemplifica un flujo de datos en un proceso sencillo.



## Figura 3. Ejemplo de MapReduce

### *Hadoop Common*

Hadoop Common Components son un conjunto de librerías y varios subproyectos de Hadoop.

Además de estos tres componentes principales de Hadoop hay otros proyectos relacionados los cuales son definidos a continuación.

### *Avro*

Es un proyecto de Apache que provee servicios de serialización y almacenamiento de datos. Cuando se guardan datos en un archivo, el esquema que define la estructura de los datos se guarda dentro del mismo; de este modo es más fácil para una aplicación leerlo posteriormente puesto que el esquema está dentro del archivo.

### *Cassandra*

Cassandra es una base de datos no relacional distribuida que utiliza un modelo de almacenamiento de <clave-valor>, desarrollado por Facebook. Permite grandes volúmenes de datos en forma distribuida y es utilizada por muchas de las empresas que utilizan Cassandra dentro de su plataforma.

### *Chukwa*

Diseñado para la colección y análisis a gran escala de datos, Chukwa es un toolkit para desplegar los resultados del análisis y monitorear el rendimiento de las aplicaciones.

### *Flume*



Tal como su nombre lo indica, su tarea principal es dirigir el flujo de datos desde una fuente hacia alguna otra localidad, en este caso hacia Hadoop. Existen tres entidades principales: *sources*, *transformers* y *sinks*. Un *source* es básicamente cualquier fuente de datos y un *sink* es el destino de una operación en específico y un *transformer* es un decorador dentro del flujo de datos que transforma esa información de una manera, como por ejemplo comprimir o descomprimir o realizar otra operación en particular sobre los mismos.

### ***HBase***

Es una base de datos columnar (column-oriented database) que ejecuta en HDFS. HBase no soporta SQL, de hecho, HBase es una base de datos relacional. Cada tabla contiene filas y columnas. HBase permite que muchos atributos de una tabla se almacenen llamándolos *familias de columnas*, de tal manera que una familia de columnas son almacenados en un solo archivo. HBase es distinto a las bases de datos relacionales orientadas a filas, ya que en las bases de datos relacionales orientadas a filas las columnas de una fila dada son almacenadas en un solo archivo. HBase utiliza HDFS en su plataforma desde Noviembre del 2009.

### ***Hive***

Es una infraestructura de data warehouse que facilita el acceso a conjuntos de datos que se encuentran almacenados en HDFS de manera distribuida. Hive tiene definido un lenguaje similar a SQL llamado HiveQL.

**Query Language(HQL), estas sentencias HQL son sep servicio de Hive y son enviadas a procesos MapReduce cluster de Hadoop.**

**El siguiente es un ejemplo en HQL para crear una tabla obtener información de la tabla utilizando Hive:**

```
CREATE TABLE Tweets (from_user STRING, userid BIGINT, tweettext STRING, retweets INT)
COMMENT "This is the Twitter feed table"
STORED AS SEQUENCEFILE;
LOAD DATA INPATH "hdfs://node/tweetdata" INTO TABLE TWEETS;
SELECT from_user, SUM(retweets)
FROM TWEETS
GROUP BY from_user;
```

## ***Jaql***

**Fue donado por IBM a la comunidad de software libre. Javascript Object Notation (JSON) es un lenguaje funcional que permite la explotación de datos en formato JSON procesar grandes volúmenes de información. Para explicar Jaql reescribe los queries de alto nivel (cuando es necesario de "bajo nivel" para distribuirlos como procesos MapReduce. Internamente el motor de Jaql transforma el query en MapReduce para reducir el tiempo de desarrollo asociado en Hadoop. Jaql posee de una infraestructura flexible para analizar datos semiestructurados como XML, archivos planos, datos relacionales, etc.**

## ***Lucene***

**Es un proyecto de Apache bastante popular para reali**

sobre textos. Lucene provee de librerías para indexar texto. Ha sido principalmente utilizado en la implementación de búsqueda (aunque hay que considerar que no tiene "crawling" ni análisis de documentos HTML ya incorporado). A nivel de arquitectura de Lucene es simple, básicamente los documentos (*document*) son divididos en campos de texto (*fields*) y se indexa sobre estos campos de texto. La indexación es el corazón de Lucene, lo que le permite realizar búsquedas rápidamente e independientemente del formato del archivo, ya sean HTML, etc.

## *Oozie*

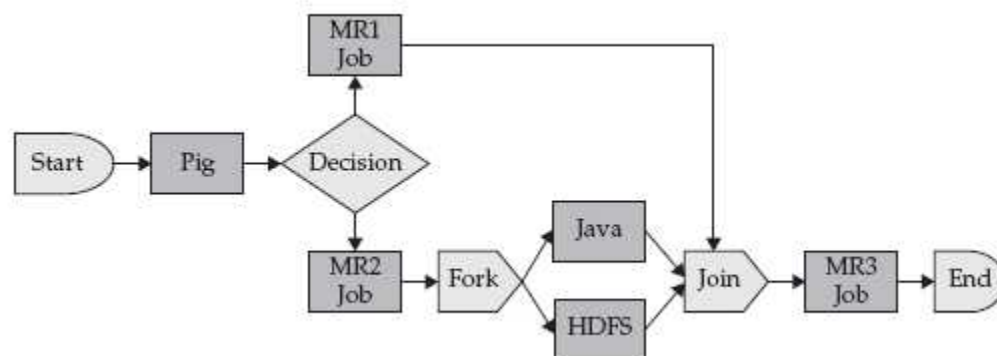
Como pudo haber notado, existen varios procesos que en distintos momentos los cuales necesitan ser orquestados para satisfacer las necesidades de tan complejo análisis de información.

Oozie es un proyecto de código abierto que simplifica la implementación y la coordinación entre cada uno de los procesos. Permite al usuario poder definir acciones y las dependencias entre dichas acciones.

Un flujo de trabajo en Oozie es definido mediante un grafo llamado *Directed Acyclical Graph (DAG)*, y es acíclico, es decir, no permite ciclos en el grafo; es decir, solo hay un punto de salida y todas las tareas y dependencias parten del punto de inicio final sin puntos de retorno. Un ejemplo de un flujo de

representa de la siguiente manera:

Figura 4. Flujo de trabajo en Oozie



### **Pig**

Inicialmente desarrollado por Yahoo para permitir a los usuarios de en analizar todos los conjuntos de datos y dedicar menos tiempo en MapReduce. Tal como su nombre lo indica al igual que cualquier ce cosa, el lenguaje *PigLatin* fue diseñado para manejar cualquier tipo ambiente de ejecución donde estos programas son ejecutados, de l relación entre la máquina virtual de Java (JVM) y una aplicación Java

### **ZooKeeper**

ZooKeeper es otro proyecto de código abierto de Apache que prove centralizada y de servicios que pueden ser utilizados por aplicacion que los procesos a través de un cluster sean serializados o sincroni Internamente en ZooKeeper una aplicación puede crear un archivo memoria en los servidores ZooKeeper llamado *znode*. Este archivo actualizado por cualquier nodo en el cluster, y cualquier nodo puede informado de los cambios ocurridos en ese *znode*; es decir, un serv configurado para "vigilar" un *znode* en particular. De este modo, las sincronizar sus procesos a través de un cluster distribuido actualiza *znode*, el cual informará al resto del cluster sobre el estatus correspo

en específico.

Como podrá observar, más allá de Hadoop, una plataforma de Big Data es un ecosistema de proyectos que en conjunto permiten simplificar, administrar y analizar grandes volúmenes de información.

## 5. Big Data y el campo de investigación

Los científicos e investigadores han analizado datos desde ya hace mucho tiempo, pero ahora representa el gran reto es la escala en la que estos son generados.

Esta explosión de "grandes datos" está transformando la manera en que se hace la investigación adquiriendo habilidades en el uso de Big Data para resolver problemas complejos relacionados con el descubrimiento científico, investigación biomédica, educación, salud, seguridad nacional, entre otros.

De entre los proyectos que se pueden mencionar donde se ha llevado a cabo una solución de Big Data se encuentran:

El *Language, Interaction and Computation Laboratory (CLIC)* en la Universidad de Trento en Italia, son un grupo de investigadores cuyo interés es estudiar la comunicación verbal y no verbal tanto con métodos computacionales como con métodos lingüísticos.

[Lineberger Comprehensive Cancer Center - Bioinformatics Group](#) utiliza Hadoop para analizar datos producidos por los investigadores de *The Cancer Research and Biotechnology Institute* para soportar las investigaciones relacionadas con el cáncer.

El [PSG College of Technology, India](#), analiza múltiples secuencias de ADN para determinar los enlaces evolutivos y predecir estructuras moleculares. El algoritmo y el paralelismo computacional de Hadoop mejora la velocidad de procesamiento de estas secuencias.

La *Universidad Distrital Francisco Jose de Caldas* utiliza Hadoop para

de investigación relacionado con el sistema de inteligencia territorial.

La *Universidad de Maryland* es una de las seis universidades que académica de cómputo en la nube de IBM/Google. Sus investigaciones en la lingüística computacional (machine translation), modelado de análisis de correo electrónico y procesamiento de imágenes.

Para más referencias en el uso de Hadoop puede dirigirse a :

<http://wiki.apache.org/hadoop/PoweredBy>

El *Instituto de Tecnología de la Universidad de Ontario (UOIT)* junto a Toronto utilizan una plataforma de big data para análisis en tiempo real (*InfoSphere Streams*), la cual permite monitorear bebés prematuros en neonatología para determinar cualquier cambio en la presión arterial, alteraciones en los registros del electrocardiograma y electroencefalograma, detectar hasta 24 horas antes aquellas condiciones que puedan ser de los recién nacidos.

Los laboratorios *Pacific Northwest National Labs (PNNL)* utilizan de *InfoSphere Streams* para analizar eventos de medidores de su red y verificar aquellas excepciones o fallas en los componentes de la red y casi de manera inmediata a los consumidores sobre el problema para administrar su consumo de energía eléctrica.[3]

La esclerosis múltiple es una enfermedad del sistema nervioso que afecta la médula espinal. La comunidad de investigación biomédica y la *Universidad de Nueva York (SUNY)* están aplicando análisis con big data para controlar la investigación, diagnóstico, tratamiento, y quizás hasta la posible cura de la esclerosis múltiple.[4]

Con la capacidad de generar toda esta información valiosa de diferentes

empresas y los gobiernos están lidiando con el problema de analizar propósitos importantes: ser capaces de detectar y responder a los cambios de una manera oportuna, y para poder utilizar las predicciones del futuro. Esta situación requiere del análisis tanto de datos en movimiento (datos en reposo (datos históricos), que son representados a diferentes y con variedades y velocidades.

## 6. Conclusiones

La naturaleza de la información hoy es diferente a la información en el pasado. La abundancia de sensores, micrófonos, cámaras, escáneres médicos, en nuestras vidas, los datos generados a partir de estos elementos se han convertido en el segmento más grande de toda la información disponible.

El uso de Big Data ha ayudado a los investigadores a descubrir cosas que se tomaron años en descubrir por sí mismos sin el uso de estas herramientas. La velocidad del análisis, es posible que el analista de datos pueda cambiar basándose en el resultado obtenido y re trabajar el procedimiento para encontrar el verdadero valor al que se está tratando de llegar.

Como se pudo notar en el presente artículo, implementar una solución de Big Data implica de la integración de diversos componentes y proyectos que forman un ecosistema necesario para analizar grandes cantidades de datos.

Sin una plataforma de Big Data se necesitaría que se desarrollara una plataforma que permita administrar cada uno de esos componentes como por ejemplo: alta conectividad, alta disponibilidad, seguridad, optimización y desempeño, monitoreo, administración de las aplicaciones, SQL y scripts personalizados. IBM cuenta con una plataforma de Big Data basada en dos productos: IBM InfoSphere BigInsights™ e IBM InfoSphere Streams™, además de IBM Watson Vivisimo, los cuales están diseñados para resolver este tipo de problemas.

herramientas están construidas para ser ejecutadas en sistemas diseñados para tratar con grandes volúmenes de información, anali estructurados como no estructurados.

Dentro de la plataforma de IBM existen más de 100 aplicaciones de trabajo que se ha realizado internamente en la empresa para casos específicas. Estos aplicativos están implementados dentro de la sol organizaciones puedan dedicar su tiempo a analizar y no a impleme

## 7. Referencias

1. Cisco, **Internet será cuatro veces más grande en 2016**, Artículo <http://www.cisco.com/web/ES/about/press/2012/2012-05-30-interr mas-grande-en-2016--informe-vini-de-cisco.html>
2. Soares Sunil, **Not Your Type? Big Data Matchmaker On Five D Explore Today**, Artículo Web <http://www.dataversity.net/not-your-data-matchmaker-on-five-data-types-you-need-to-explore-today/>
3. Clegg Dai, **Big Data: The Data Velocity Discussion**, Artículo We <http://thinking.netezza.com/blog/big-data-data-velocity-discussion>
4. Kobielus James, **Big Data Analytics Helps Researchers Drill D Sclerosis**, Artículo Web <http://thinking.netezza.com/blog/big-data researchers-drill-deeper-multiple-sclerosis>
5. Aprenda más acerca de Apache Hadoop en <http://hadoop.apache>
6. Zikopolous Paul, Deroos Dirk, Deutsch Tom, Lapis George, **Unde Analytics for Enterprise Class Hadoop and Streaming Data**, M
7. Foster Kevin, Nathan Senthil, Rajan Deepak, Ballard Chuck, **IBM Assembling Continuous Insight in the Information Revolutior**







## **Bluemix**

Guías rápidas de inicio y demos de la plataforma abierta en la Nube de IBM.



## **Súbete a la Nube de IBM**

Abre el potencial del cómputo en la nube con productos y servicios de IBM.



## **Arquitectura de Nube Abierta de IBM**

Al cambiar cómo se manejan los negocios y la sociedad, la computación en nube está abriendo gigantescas avenidas de innovaciones.